



Eur päisches Patentamt  
European Patent Office  
Office européen des brevets



(11) EP 0 452 413 B1

(12) EUROPEAN PATENT SPECIFICATION

- (45) Date of publication and mention of the grant of the patent: 12.04.2000 Bulletin 2000/15
- (21) Application number: 90902453.1
- (22) Date of filing: 05.01.1990
- (51) Int. Cl.<sup>7</sup>: C12N 15/00, C12Q 1/08
- (86) International application number: PCT/US90/00024
- (87) International publication number: WO 90/07862 (26.07.1990 Gazette 1990/17)

(54) GENERATION AND SELECTION OF NOVEL DNA-BINDING PROTEINS AND POLYPEPTIDES  
HERSTELLUNG UND SELEKTION VON DNA-BINDUNGSPROTEINEN UND POLYPEPTIDEN  
PRODUCTION ET SELECTION DE PROTEINES ET DE POLYPEPTIDES NOUVEAUX DE LIAISON A L'ADN

- (84) Designated Contracting States:  
AT BE CH DE DK ES FR GB IT LI LU NL SE
- (30) Priority: 06.01.1989 US 293980
- (43) Date of publication of application: 23.10.1991 Bulletin 1991/43
- (73) Proprietor: Dyax Corp.  
Cambridge, MA 02139 (US)
- (72) Inventors:  
• LADNER, Robert, Charles  
Ijamsville, MD 21754 (US)  
• GUTERMAN, Sonia, Kosow  
Belmont, MA 02178 (US)  
• KENT, Rachel, Baribault  
Wilmington, MA 01887 (US)  
• LEY, Arthur, Charles  
Newton, MA 02165 (US)
- (74) Representative:  
Plougmann, Ole et al  
Plougmann, Vingtoft & Partners,  
Sankt Annae Plads 11,  
P.O. Box 3007  
1021 Copenhagen K (DK)

- (56) References cited:  
WO-A-88/06601
- SCIENCE, vol. 242, 14 October 1988, AAAS, WASHINGTON, DC, US; pages 240 - 245; S. BASS ET AL.: 'Mutant trp repressor with new DNA-binding specificities'
  - PROC. NATL. ACAD SCI. vol. 85, August 1988, NATL. ACAD SCI., WASHINGTON, DC, US; pages 5834 - 5838; M. HOLLIS ET AL.: 'A receptor heterodimer binds to a chimeric operator'
  - J. CELL. BIOCHEM. SUPPL. 0 (13 PART D), 18TH UCLA SYMPOSIA ON MOLECULAR AND CELLULAR BIOL., PARK CITY, UTAH, US; 1-7 APRIL 1989, ABSTRACT NO. M334; page 307; D. PIOLI ET AL.: 'Control of plant gene expression using wild-type and altered-specificity bacterial repressor molecules'
  - Biochemistry, Volume 28, No. 23, Issued 1989, HUGHES et al., "Purification and Characterization of a protein from Hela cells that binds with high affinity to the estrogen response element, GGTCAGCGTGACC" See pages 9137-9142.
  - Proceedings of the National Academy of Sciences USA Volume 85 issued December 1988 GAYNOR et al. "Repeated B motifs in the human immunodeficiency virus type I long terminal repeat enhancer region do not exhibit cooperative factor binding" See pages 9406-9410.
  - Journal of Virology Volume 63(6) Issued June 1989 HARRICH et al. "Role of SPI-binding domains in the in vivo transcriptional regulation of the human immunodeficiency virus type 1 long terminal repeat" See pages 2585-2591.

Note: Within nine months from the publication of the mention of the grant of the European patent, any person may give notice to the European Patent Office of opposition to the European patent granted. Notice of opposition shall be filed in a written reasoned statement. It shall not be deemed to have been filed until the opposition fee has been paid. (Art. 99(1) European Patent Convention).

EP 0 452 413 B1

- 
- |  |  |
|--|--|
| <ul style="list-style-type: none"><li>• Biological Abstracts Volume 87(4) Issued 15 February 1989 SWEDER et al. "Purification and characterization of proteins that bind to yeast ARSs" See ref. 37597.</li><li>• Biological Abstracts Volume 85(9) Issued 01 May 1988 HAUBER et al. "Mutational analysis of the trans-activation responsive region of the human immunodeficiency virus type long terminal repeat" See reference 89626.</li><li>• Life Sciences Collection Accession Number 82001774133 Issued 1988 HENDERSON "Isolation and characterization of a novel protein (X-ORF) product from ISIV and HIV-2". See Abstract.</li></ul> | <ul style="list-style-type: none"><li>• Life Sciences Collection Accession Number 82001901798 Issued 1988 "Purification of the human immunodeficiency virus type 1 enhancer and TAR binding proteins EBP-1 and UBP-1". See Abstract.</li><li>• Genes and Development (1989) 3: 185-197</li><li>• Proc. Natl. Acad. Sci. USA (1989) 86: 3689-3693</li><li>• Genetics (1986) 114: 1-14</li></ul> |
|--|--|
-

## Description

## BACKGROUND OF THE INVENTION

5 Field of the Invention

[0001] This invention relates to development of novel DNA-binding proteins and polypeptides by an iterative process of mutation, expression, selection, and amplification. The ability to create novel DNA-binding proteins will have far-reaching applications, including, but not limited to, use in: a) treating viral diseases, b) treating genetic diseases, c) preparation of novel biochemical reagents, and d) biotechnology to regulate gene expression in cell cultures. Several workers have shown that repressors derived from bacteria function when expressed in eukaryotic cells (BREN84, FIGG88, BROW87, HUMC87, HUMC88), but none have shown how to generate proteins that bind sequence-specifically to a predetermined DNA sequence. For reviews of transcriptional control in eukaryotic cells, see STRU87, JONE87, and MANI87. The present application deals only with sequence-specific DNA-binding proteins, abbreviated DBP.

[0002] Proteins, particularly repressors, having affinity for specific sites on DNA modulate transcription of genes. The best known are a group of proteins primarily studied in prokaryotes that contain the structural motif alpha-helix-turn-alpha-helix (H-T-H) (SAUE82, PABO84). These proteins bind as dimers or tetramers to DNA at specific operator sequences that have approximately palindromic sequences. Contacts made by two adjacent alpha helices of each monomer in and around two sites in the major groove of B-form DNA are a major feature in the DNA-protein interface. This group of proteins includes phage repressor and Cro proteins, bacterial metabolic repressors such as GalP, LacI, LexA, and TrpR, bacterial activator protein CAP and activator/repressor AraC, bacterial transposon and plasmid TetR proteins (PABO84), the yeast mating type regulators MATa1 and MATalpha2 (MILL85) and eukaryotic homeo box proteins (EVAN88).

[0003] Interactions between dimeric repressors and approximately palindromic operators have usually been discussed in the literature with attention focused on one half of the operator with the tacit or explicit assumption that identical interactions occur in each half of the complex. Departures from palindromic symmetry allow proteins to distinguish among multiple related operators (SADL83, SIMO84). One must view the DNA-protein interface as a whole. The emphasis in the literature on dyad symmetry is a barrier to determining the requirements for general novel recognition of DNA by proteins.

[0004] The equilibrium geometry and flexibility of DNA are determined by the sequence; see *inter alia* HOGA87, GART88, and ULAN87. The interactions of ionic, polar, and hydrophobic groups on the DNA with solvent molecules and ions make detailed predictions of DNA conformation and binding properties very difficult; cf. OHLE85, ULAN87, and OTWI88.

[0005] Matthews (MATT88), commenting on the current collection of protein-DNA structures, concludes that: a) different H-T-H DBPs use their recognition helices differently, b) there is no simple code that relates particular base pairs to particular amino acids at specific locations in the DBP, and c) "full appreciation of the complexity and individuality of each complex will be discouraging to anyone hoping to find simple answers to the recognition problem." Schleif (SCHL88) has characterized the study of DNA-binding proteins as a field still in its infancy and emphasizes the difficulties of designing proteins that bind predetermined sequences.

[0006] Prokaryotic repressors exist that are unrelated to H-T-H binding proteins. Some of these bind to approximate palindromic sequences (e.g. *Salmonella typhimurium* phage P22 Mnt protein (VERS87a) and *E. coli* TyrR repressor protein (DEFE86)). Others bind to operator sequences that are partially symmetric (*S. typhimurium* phage P22 Arc protein, VERS87b; *E. coli* Fur protein, DELO87; plasmid R6K pi protein, FILU85) or non-symmetric (phage Mu repressor, KRAU86).

[0007] Genetics has enabled extensive analysis of prokaryotic DNA-binding proteins and their specific nucleic acid recognition sequences. It is not yet possible, however, to design a protein to bind strongly and specifically to an arbitrary DNA sequence. As taught by the present invention it is, nonetheless, possible to postulate a family of potential DBP mutants and identify one having the desired specificity by other means.

[0008] Genetic studies of the DNA-binding proteins show that mutations in protein sequence that result in decrease of protein function fall into two overlapping classes: 1) those that destabilize the global protein structure or folding and 2) those that specifically alter the binding properties. The first class illuminates the general problem of protein folding and stability, while the second defines the interactions involved in the formation and stabilization of the protein-DNA complex. Mutations in the operator yield additional information.

[0009] Positions 84 to 91 in helix 5 of  $\lambda$  repressor have been subjected to extensive amino acid substitutions (REID88). Two or three positions were varied simultaneously through all twenty amino acids and those combinations giving normal function were selected. The authors neither discuss optimization of the number or positions of residues to vary to obtain any particular functionality, nor did they attempt to obtain proteins having alternate dimerization or rec-

ognition functions.

[0010] Pakula *et al.* (PAKU86) have randomly mutagenized  $\lambda$  Cro. They sought and found non-functional mutants but did not seek or find proteins having novel DNA-binding properties, nor did they suggest how to select such proteins.

[0011] Sequence-independent DNA-protein interactions are thought to occur via electrostatic interactions between the backbone of the DNA and charged or polar groups of the protein (ANDE87, LEWI83, and TAKE85). Sequence-specific interactions involve H-bonding, nonpolar, or van der Waals contacts between exposed side groups or groups of the polypeptide main chain and base pair edges exposed in the major and minor grooves of the DNA.

[0012] Mutations that alter residues involved in specific binding interactions with DNA have been identified in prokaryotic DBPs, including  $\lambda$ , 434, and P22 repressor and Cro proteins, P22 Arc and Mnt, and *E. coli* *trp* and *lac* repressors and CAP. These mutations occur in residues that are exposed to solvent in the free protein but buried in the protein-DNA complex.

[0013] A few cases have been reported (BASS88, YOUN83, VERS85a, CARU87, WHAR85b, WHAR87, EBRI84, and SPIR88) in which a change in one or a few residues in a DNA-binding protein not only abolishes binding by the protein to the wild-type operator but also confers strong binding to a different operator. In all the cited publications, alteration of binding specificity has been accomplished by using symmetrically-located pairs of alterations in the operator sites. Single, asymmetric changes or multiple changes asymmetrically located in either the binding protein or its operator were not considered. In "helix swap" experiments (WHAR84, WHAR85b, WHAR85a, SPIR88, BUSH88, PABO84), multiple mutations are introduced into the DNA-binding recognition helix of H-T-H proteins with the goal of changing the operator specificity of one known DBP to that of a different known DBP.

[0014] An extension of the "helix swap" experiments uses a mixture of 434 repressor and 434R(alpha3(P22R)) (HOLL88). This mixture recognizes and binds *in vitro* with high affinity to a 16 bp chimeric operator comprising a 434 half-site and a P22 half-site, indicating that active heterodimers are formed. The authors did not extend the results to intracellular repression, nor did they perform mutagenesis of the repressors and selection of cells to create novel recognition patterns.

[0015] Two approaches have been developed to create novel proteins through reverse genetics. In one approach, dubbed "protein surgery" (DILL87), a substitution is introduced at a single protein residue. This approach has been used to determine the effects on structure and function of specific substitutions in trypsin (CRAI85, RAOS87, BASH87).

[0016] The other approach has been to generate a variety of mutants at many loci within the cloned gene, the "gene-directed random mutagenesis" method. The specific location and nature of the change or changes are determined *post hoc* by DNA sequencing. If loss of a wild-type function confers a cellular phenotype, one screens colonies for mutations; see, *cf.* PAKU86. This approach is limited by the number of colonies that can be examined. An additional important limitation is that many desirable protein alterations require multiple amino acid substitutions and thus are not accessible through single base changes or even through all possible amino acid substitutions at any one residue.

[0017] The objective in both these approaches has been, however, to analyze the effects of a variety of point mutations, so that rules governing such substitutions could be developed (ULME83). Progress has been hampered by the efforts involved in using either method (ROBE86).

[0018] Oliphant *et al.* (OLIP86) and Oliphant and Struhl (OLIP87) have demonstrated ligation and cloning of highly degenerate oligonucleotides and have applied saturation mutagenesis to the study of promoter sequence and function. They have suggested that similar methods could be used to study genetic expression of protein coding regions of genes, but they do not say how one should: a) choose protein residues to vary, or b) select or screen mutants with desirable properties.

[0019] Ward *et al.* (WARD86) have engineered heterodimers from homodimers of tyrosyl-tRNA synthetase. Methods of converting homodimeric DBPs into heterodimeric DBPs are disclosed in the present invention. Methods of deriving single-polypeptide pseudo-dimeric DBPs from homodimeric DBPs are disclosed in the examples of the present invention.

[0020] Benson *et al.* (BENS86) have developed a scheme to detect genes for sequence-specific DNA-binding proteins. They do not consider non-symmetric target DNA sequences nor do they suggest mutagenesis to generate novel DNA-binding properties. Their method is presented as a method to detect genes for naturally occurring DNA-binding proteins. Because the selective system is lytic growth of phage, low levels of repression can not be detected. Selective chemicals, as disclosed in the present application, on the other hand, can be finely modulated so that low level repression is detectable.

[0021] Elledge and Davis (ELLE89a) and Elledge *et al.* (ELLE89b) have used an occluded *aadA* gene in a selection for cells expressing eukaryotic DBPs. The supposed recognition sequence of the sought DBP is incorporated into the strong promoter that occludes *aadA* on a low-copy number plasmid. Their system is presented as a tool for cloning pre-existing DBPs and there is no mention of variegation of the gene that encodes the potential DBP. Furthermore, there is no discussion of the symmetry of the target sequence or of the symmetry of the DBP.

[0022] Ladner and Bird, WO88/06601 suggest strategies for the preparation of asymmetric repressors. In one embodiment, a gene is constructed that encodes, as a single polypeptide chain, the two DNA-binding domains of a nat-

usually-occurring dimeric repressor, joined by a polypeptide linker that holds the two binding domains in the necessary spatial relationship for binding to an operator. While they prefer to design the linker based on protein structural data (cf. Ladner, U.S. Patent 4,704,692) they state that uncertainties in the design of the linker may be resolved by generating a family of synthetic genes, differing in the linker-encoding subsequence, and selecting in vivo for a gene encoding the desired pseudo-dimer. Ladner and Bird do not consider the background of false positives that would arise if the two-domain polypeptides dimerize to form pseudo-tetramers.

[0023] The binding of lambdoid repressors, Cro and CI repressor, is taken, in WO88/06601, as canonical even though other DBPs were known having operators of different lengths. WO88/06601 maintains that the 17 bp lambdoid operators can be divided into three regions: a) a left arm of five bases, b) a central region of seven bases, and c) a right arm of five bases. Several other DBPs are known for which this division is inappropriate. Further, WO88/06601 states that the sequence and composition of the central region, in which edges of bases are not contacted by the DBP, are immaterial. There is direct evidence for 434 repressor (K0UD87, K0UD88) that the sequence and composition of the central region strongly influences binding of 434 repressor.

[0024] Once a pseudo-dimer is obtained, they then obtain an asymmetric pseudo-dimer by the following technique. First, the user of WO88/06601 is directed to construct a family of hybrid operators in which the sequence of the left and right arms are specified; no specification is given for the central seven bases. In each member of the family, the left arm contains the same sequence as the wild-type operator left arm while the right arm 5-mer is systematically varied through all 1024 possibilities. Similarly, in the gene encoding the pseudodimer, the codons for one recognition helix have the wild-type sequence while the codons coding for the other recognition helix are highly varied. The variegated pseudodimer genes are expressed in bacterial cells, wherein the hybrid operators are positioned to repress a single highly deleterious gene. Thus, it is supposed that one can identify a recognition helix for each possible 5-mer right arm of the operator by in vivo selection; the correspondences between 5-mer right arms and sequences of recognition helices are compiled into a dictionary. The consequences of mutations or deletions in the deleterious genes are not considered. WO88/06601 suggests that successful constructions may be very rare, e.g. one in  $10^6$ , but ignore other genetic events of similar or greater frequency.

[0025] To obtain a repressor for an arbitrary 17-mer operator, the user of WO88/06601:

a) finds the 5-mer sequence of the left arm in the dictionary and uses the corresponding recognition helix sequence in the first DNA-binding domain of the pseudodimer,

b) ignores the sequence and composition of the next seven bases, and

c) finds the 5-mer sequence of the right arm in the dictionary and uses the corresponding recognition helix sequence in the second DNA-binding domain of the pseudodimer.

[0026] WO88/06601 also envisions means for producing a heterodimeric repressor. A plasmid is provided that carries genes encoding two different repressors. A population of such plasmids is generated in which some codons are varied in each gene. WO88/06601 instructs the user to introduce very high levels of variegation without regard to the number of independent transformants that can be produced. WO88/06601 also instructs the user to introduce variegation at widely separated sites in the gene, though there is no teaching concerning ways to simultaneously introduce high levels of variegation at widely separated sites in the gene or concerning maintenance of diversity without selective pressure, as would be needed if the variegation were introduced stepwise. WO88/06601 teaches that codons thought to be involved in the protein-protein interface should be preferentially mutated to generate heterodimers. Cells transformed with this population of plasmids will produce both the desired heterodimer and the two "wild-type" homodimers. WO88/06601 advises that one select for production of the heterodimer by providing a highly deleterious gene controlled by a hybrid operator, and beneficial genes controlled by the wild-type operators. The fastest growing cells, it is taught, will be those that produce a great deal of the heterodimer (which blocks expression of the deleterious gene) and little of the homodimers (so that the beneficial genes are more fully expressed). There is no consideration of mutations or deletions in the deleterious gene or in the wild-type operators; such mutations will produce a background of fast-growing cells that do not contain the desired heterodimers.

## SUMMARY OF THE INVENTION

[0027] This invention relates to the development of novel proteins or polypeptides that preferentially bind to a specific subsequence of double-stranded DNA (the "target") using a novel scheme for in vivo selection of mutant proteins exhibiting the desired binding specificities.

[0028] The invention relates in particular to a selection vector for selecting recipient cells transformed by such vector that express a protein or polypeptide that binds specifically to a predetermined target DNA sequence borne by said

vector, which comprises a first and a second operon, each comprising at least one expressible gene, the genes of said first and second operon being different, a copy of the target DNA sequence being included in each operon and positioned therein so that the recipient cells enjoy a selective advantage, other than resistance to lytic growth of phage, if they express a protein or polypeptide which binds to said copies of the target DNA sequence.

[0029] The novel binding proteins or polypeptides may be obtained by mutating a gene encoding on expression: 1) a known DNA-binding protein within the subsequence encoding a known DNA-binding domain, 2) a protein that, while not possessing a known DNA-binding activity, possesses a secondary or higher order structure that lends itself to binding activity (clefs, grooves, helices, etc.), 3) a known DNA-binding protein but not in the subsequence known to cause the binding, or 4) a polypeptide having no known 3D structure of its own.

[0030] This application uses the term "variegated DNA" to refer to a population of molecules that have the same base sequence through most of their length, but that vary at a number of defined loci. Using standard genetic engineering techniques, variegated DNA can be introduced into a plasmid so that it constitutes part of a gene (OLIP86, OLIP87, CHEN88, AUSU87, REID88). When plasmids containing variegated DNA are used to transform bacteria, each cell makes a version of the original protein. Each colony of bacteria produces a different version from most other colonies. If the variegations of the DNA are concentrated at loci that code on expression for residues known to be on the surface of the protein or in loops, a population of genes will be generated that code on expression for a population of proteins, many members of which will fold into roughly the same 3D structure as the parental protein. Most often we generate mutations that are concentrated within codons for residues thought to make contact with the DNA. Secondly, we introduce mutations into codons specifying residues that are not directly involved in DNA contact but that affect the position or dynamics of residues that do contact the DNA.

[0031] In general, a variegated population of DNA molecules, each of which encodes one of a large (e.g.  $10^7$ ) number of distinct potential target-binding proteins, is used to transform a cell culture. The cells of this cell culture are engineered with binding marker genetic elements so that, under selective conditions, the cell thrives only if the expressed potential target-binding protein in fact binds to the target subsequence preventing transcription of these binding marker genetic elements. (Typically, binding of a successful target-binding protein to the target subsequence blocks expression of a gene product that is deleterious under selective conditions. Alternatively, binding of a successful target-binding protein can inactivate a strong promoter that otherwise occludes transcription of a beneficial gene.) The mutant cells are directed to express the potential target-binding proteins and the selective conditions are applied. Cells expressing proteins binding successfully to the target are thus identified by *in vivo* selection. If the binding characteristics are not fully satisfactory, the amino acid sequences of the best binding proteins are determined (usually by sequencing the corresponding genes), a new population of DNA molecules is synthesized that encode variegated forms of the best binding proteins of the last cull, mutant cells are prepared, the new population of potential DNA-binding proteins is expressed, and the best proteins are once again identified by the superior growth of the corresponding transformants under selective conditions. The process is repeated until a protein or polypeptide with the desired binding characteristics is obtained. Its corresponding gene may then be moved to a suitable expression system.

[0032] In the simplest form of this invention, the mutant cells are provided with a selectable genetic element, the transcription of which is deleterious to the survival or growth of the cell. The selectable genetic element either is a promoter or is operably linked to a promoter regulating the expression of the gene. The promoter, or other non-coding region of the genetic element (for example, an intron), has been modified to include the desired target subsequence in a position where it will not interfere with transcription of the selectable gene unless a protein binds to that target subsequence. Each mutant cell is also provided with a gene encoding on expression a potential DNA-binding protein, operably linked to a promoter that is preferably regulated by a chemical inducer. When this gene is expressed, the potential DNA-binding protein has the opportunity to bind to the target and thereby protect the cell from the selective conditions under which the product of the binding marker gene would otherwise harm the cell.

[0033] In addition to the desired outcome of these *in vivo* selections, there exist a number of possible genetic events that allow the cells to escape the selection, producing artifacts and inefficiency by allowing the growth of colonies that do not express the desired sequence-specific DNA-binding proteins. Examples of mechanisms, other than the desired outcome, that lead to cell survival under the selective conditions include: a) a point mutation or a deletion in the selectable gene eliminates expression or function of the selectable gene product; b) a host chromosomal mutation compensates for or suppresses function of the selectable gene product; c) the introduced potential DNA-binding protein binds to a DNA subsequence other than the chosen target subsequence and blocks expression of the selectable gene; d) the introduced potential DNA-binding protein binds to and inactivates the gene product of the selective gene; and e) a DNA-binding protein endogenous to the host mutates so that it binds to the selectable gene and blocks expression of the selectable gene.

[0034] This invention relates, in particular, to the design of a vector that confers upon the host cells the desired conditional sensitivity to the selection conditions in such a manner as to greatly reduce the likelihood of false positives and artifactual colonies.

[0035] First, at least two selectable genes that are functionally unrelated are used to reduce the risk that a single

point mutation in the vector (or in the host chromosome) will destroy the sensitivity of the cell to the selective conditions, since it will eliminate only one of the two (or more) deleterious phenotypes. Similarly, a single introduced gene for a potential DNA-binding protein that binds to and inactivates the gene product of one selectable gene will not bind and inactivate the gene product of the other selectable gene. The likelihood that point mutations will occur in both selectable genes or that two host chromosomal mutations will spontaneously arise that suppress the effects of two genes is the product of each single individual probabilities of the necessary event, and thus is extremely low.

[0036] The DNA sequences of the two or more selectable genes preferably should not have long segments of identity: a) to avoid isolation of a DBP that binds these identical regions instead of the intended target sequence, and b) to reduce the likelihood of genetic recombination. The degeneracy of the genetic code allows us to avoid exact identity of more than a few, e.g., 10, bases.

[0037] Second, the selectable genes are placed on the vector in alternation with genetic elements that are essential to plasmid maintenance. Thus, a single deletion event, even of thousands of bases, cannot eliminate both selectable genes without also eliminating vital genetic elements. Alternatively, the selectable genes are placed in the bacterial chromosome. Spontaneous deletions from the chromosome are rare.

[0038] Third, different promoters are associated with each of the selectable genes. This ensures that the selection does not isolate cells harboring genes encoding on expression novel DNA-binding proteins that bind specifically to subsequences that are part of the promoter but not the chosen target subsequence. Each cell expresses only one or a few introduced potential DNA-binding proteins (multiple potential DNA-binding proteins could arise if one cell is transformed by two or more variegated plasmids). The probability that two such proteins will occur in one cell and that one will bind to the promoter of the first selectable gene and that the second will bind to the different promoter of the second selectable gene is very small.

[0039] Fourth, the selectable binding marker genes may be placed on a vector different from the vector that carries the potential *dbp* gene. DNA manipulations that introduce variegation into the potential *dbp* gene can cause mutations in the vector remote from the site of the intended mutations. Thus, we may place the selectable binding marker genes in the bacterial chromosome or on a separate plasmid that is compatible with the *dbp* vector.

[0040] Finally, the same promoter is used to initiate transcription of two genes: a) one of the deleterious selectable binding marker genes, and b) a beneficial or essential gene also borne on the plasmid and used to select for uptake and maintenance of the plasmid (e.g., an antibiotic resistance gene, such as *bla*). In the case of the beneficial or essential gene, however, there is no instance of the predetermined target DNA subsequence associated with the promoter. Thus, if a DNA-binding protein binds to a subsequence of the promoter other than the predetermined target DNA subsequence, it will frustrate expression of the beneficial or essential one. If desired, more than one such beneficial or essential gene may be provided. In that event, copies of promoter A may be operably linked to both deleterious gene A' (with an instance of the target) and beneficial gene A'' (without an instance of the target), while copies of promoter B are operably linked to both deleterious gene B' (with target) and beneficial gene B'' (without target).

[0041] The selection system described above is a powerful tool that eliminates most of the artifacts associated with selections based on cloning vectors that use a single selectable gene or that have all selectable genes in a contiguous region of the plasmid. While this invention embraces using the aforementioned elements of a selection system singly or in partial combination, most preferably all are employed.

[0042] In one embodiment, the invention relates to a cell culture comprising a plurality of cells, each cell bearing:

i) a gene coding on expression for a potential DNA-binding protein or polypeptide, where such protein or polypeptide is not the same for all such cells, but rather varies at a limited number of amino acid positions; and

ii) at least two independent operons, each comprising at least one binding marker gene coding on expression for a product conditionally deleterious to the survival or reproduction of such cells, the promoter of each said binding marker gene containing a predetermined target DNA subsequence so positioned that, if said target DNA subsequence is bound by a DNA-binding protein or polypeptide, said conditionally deleterious product is not expressed in functional form.

[0043] Abolition of function is much easier than engineering of novel function. Reverse selection can isolate cells that: a) express no DBP, b) express unstable proteins descendant from a parental DBP, c) express a protein descendant from a parental DBP having very nearly the same 3D structure as the parental DBP, but lacking the functionality of the parent. We are interested in this third class. It is difficult, however, to distinguish among these classes genetically. Therefore, when using reverse selection, we carefully choose sites to mutate the protein (so as to minimize the chances of destroying tertiary structure) and we introduce a lower level of variegation than in forward selection. We must verify biochemically that a stable, folded protein is produced by the isolated cells.

[0044] Another concept of the present invention is the use of a polypeptide, rather than a protein, to preferentially bind DNA. This polypeptide, instead of binding the DNA molecule as a preformed molecule having shape complement-

tary to DNA, will wind about the DNA molecule in the major or minor groove. Such a polypeptide has the advantages that: a) it is smaller than a protein having equivalent recognizing ability and may be easier to introduce into cells, and b) it may serve as a model for creation of other compounds that bind DNA sequence-specifically.

[0045] In a preferred embodiment, transcription of the DNA that codes on expression for potential-DNA-binding proteins or polypeptides is regulated by addition of chemical inducer to the cell culture, such as isopropylthiogalactoside (IPTG). Other regulatable promoters having different inducers or other means of regulation are also appropriate.

[0046] The invention encompasses the design and synthesis of variegated DNA encoding on expression a collection of closely related potential DNA-binding proteins or polypeptides characterized by constant and variable regions, said proteins or polypeptides being designed with a view toward obtaining a protein or polypeptide that binds a prede-

termined target DNA subsequence.

[0047] For the purposes of this invention, the term "potential DNA-binding polypeptide" refers to a polypeptide encoded by one species of DNA molecule in a population of variegated DNA wherein the region of variation appears in one or more subsequences encoding one or more segments of the polypeptide having the potential of serving as a DNA-binding domain for the target DNA sequence or having the potential to alter the position or dynamics of protein residues that contact the DNA. A "potential DNA-binding protein" (potential-DBP) may comprise one or more potential DNA-binding polypeptides. Potential-DBPs comprising two or more polypeptide chains may be homologous aggregates (e.g. A<sub>2</sub>) or heterologous aggregates (e.g. AB).

[0048] From time to time, it may be helpful to speak of the "parental sequence" of the variegated DNA. When the novel DNA-binding domain sought is a homolog of a known DNA-binding domain, the parental sequence is the sequence that encodes the known DNA-binding domain. The variegated DNA is identical with this parental sequence at most loci, but will diverge from it at chosen loci. When a potential DNA-binding domain is designed from first principles, the parental sequence is a sequence that encodes the amino acid sequence that has been predicted to form the desired DNA-binding domain, and the variegated DNA is a population of "daughter DNAs" that are related to that parent by a high degree of sequence similarity.

[0049] The fundamental principle of the invention is one of forced evolution. The efficiency of the forced evolution is greatly enhanced by careful choice of which residues are to be varied. The 3D structure of the potential DNA-binding domain and the 3D structure of the target DNA sequence are key determinants in this choice. First a set of residues that can either simultaneously contact the target DNA sequence or that can affect the orientation or flexibility of residues that can touch the target is identified. Then all or some of the codons encoding these residues are varied simultaneously to produce a variegated population of DNA. The variegated population of DNA is introduced into cells so that a variegated population of cells producing various potential-DBPs is obtained.

[0050] The highly variegated population of cells containing genes encoding potential-DBPs is selected for cells containing genes that express proteins that bind to the target DNA sequence ("successful DNA-binding proteins"). After one or more rounds of such selection, one or more of the chosen genes are examined and sequenced. If desired, new loci of variation are chosen. The selected daughter genes of one generation then become the parental sequences for the next generation of variegated DNA (vgDNA).

[0051] DNA-binding proteins (DBPs) that bind specifically to viral DNA so that transcription is blocked will be useful in treating viral diseases, either by introducing DBPs into cells or by introducing the gene coding on expression for the DBP into cells and causing the gene to be expressed. In order to develop such DBPs, we need use only the nucleotide sequence of the viral genes to be repressed. Once a DBP is developed, it is tested against virus in vivo. Use of several independently-acting DBPs that all bind to one gene allow us to: a) repress the gene despite possible variation in the sequence, and b) to focus repression on the target gene while distributing side effects over the entire genome of the host cell. Animals, plants, fungi, and microbes can be genetically made intracellularly immune to viruses by introducing, into the germ line, genes that code on expression for DBPs that bind DNA sequences found in viruses that infect the animal (including human), plant, fungus, or microbe to be protected.

[0052] Sequence-specific DBPs may also be used to treat autoimmune and genetic disease either by repressing noxious genes or by causing expression of beneficial genes.

[0053] Some naturally-occurring DBPs bind sequence-specifically to DNA only in the presence or absence of specific effector molecules. For example, Lac repressor does not bind the lac operator in the presence of lactose or isopropylthiogalactoside (IPTG); Trp repressor binds DNA only in the presence of tryptophan or certain analogues of tryptophan. The method of the present invention can be used to select mutants of such DBPs that a) recognize a different cognate DNA sequence, or b) recognize a different effector molecule. These alterations would be useful because: a) known inducible or de-repressible DBPs allows us to use the novel DBP without affecting existing metabolic pathways. Having novel effectors allows us to induce or de-repress the regulated gene without altering the state of genes that are controlled by the natural effectors. In addition, temperature-sensitive DBPs could be made which would allow us to control gene expression in the same way that  $\lambda$  cI857 and P<sub>R</sub> and P<sub>L</sub> are used.

[0054] Conferring novel DNA-recognition properties on proteins will allow development of novel restriction enzymes that recognize more base pairs and therefore cut DNA less frequently. For example, the methods of the present inven-



tion will be useful in developing a derivative of EcoRI (recognition GAATTC) that recognizes and cleaves a longer recognition site, such as TGAATTC. Proteins that recognize specific DNA sequences may also be used to block the action of known restriction enzymes at some subset of the recognition sites of the known enzyme, thereby conferring greater specificity on that enzyme. Other DNA-binding enzymes may also be obtained by the methods described herein.

[0055] The methods of the present invention are primarily designed to select from a highly variegated population those cells that contain genes that code on expression for proteins that bind sequence-specifically to predetermined DNA sequences. The genetic constructions employed can also be used as an assay for putative DBPs that are obtained in other ways.

#### BRIEF DESCRIPTION OF THE DRAWINGS

[0056]

- Figure 1 Schematic of protein bound to DNA.
- Figure 2 Schematic of evolution of a binding protein.
- Figure 3 Plasmid pKK175-6.
- Figure 4 Plasmid pAA3H.
- Figure 5 Summary of construction of pEP1009.
- Figure 6 Plasmid pEP1001.
- Figure 7 Plasmid pEP1002.
- Figure 8 Plasmid pEP1003.
- Figure 9 Plasmid pEP1004.
- Figure 10 Plasmid pEP1005.
- Figure 11 Plasmid pEP1007.
- Figure 12 Plasmid pEP1009.

#### DETAILED DESCRIPTION OF THE INVENTION AND ITS PREFERRED EMBODIMENTS

##### Abbreviations :

[0057] The following abbreviations will be used throughout the present invention:

| Abbreviation                        | Meaning                                   |
|-------------------------------------|---|
| DBP                                 | DNA-binding protein                       |
| <u>idbp</u>                         | A gene encoding the initial DBP           |
| <u>pdbp</u>                         | A gene encoding a potential-DBP           |
| vgDNA                               | variegated DNA                            |
| dsDNA                               | double-stranded DNA                       |
| ssDNA                               | single-stranded DNA                       |
| Tc <sup>R</sup> , Tc <sup>S</sup>   | Tetracycline resistance or sensitivity    |
| Gal <sup>R</sup> , Gal <sup>S</sup> | Galactose resistance or sensitivity       |
| Gal <sup>+</sup> , Gal <sup>-</sup> | Ability or inability to utilize galactose |
| Fus <sup>R</sup> , Fus <sup>S</sup> | Fusaric acid resistance or sensitivity    |
| Km <sup>R</sup> , Km <sup>S</sup>   | Kanamycin resistance or sensitivity       |
| Ap <sup>R</sup> , Ap <sup>S</sup>   | Ampicillin resistance or sensitivity      |

##### Terminology

[0058] A domain of a protein that is required for the protein to specifically bind a chosen DNA target subsequence.

is referred to herein as a "DNA-binding domain". A protein may comprise one or more domains, each composed of one or more polypeptide chains. A protein that binds a DNA sequence specifically is denoted as a "DNA-binding protein". In one embodiment of the present invention, a preliminary operation is performed to obtain a stable protein, denoted as an "initial DBP", that binds one specific DNA sequence. The present invention is concerned with the expression of numerous, diverse, variant "potential-DBPs", all related to a "parental potential-DBP" such as a known DNA-binding protein, and with selection and amplification of the genes encoding the most successful mutant potential-DBPs. An initial DBP is chosen as parental potential-DBP for the first round of variegation. Selection isolates one or more "successful DBPs". A successful DBP from one round of variegation and selection is chosen to be the parental DBP to the next round. The invention is not, however, limited to proteins with a single DNA binding domain since the method may be applied to any or all of the DNA binding domains of the protein, sequentially or simultaneously.

[0059] Amino acids are indicated by the single-letter code, AUSU87, Appendix A.

[0060] Symbols that represent ambiguous DNA are: T, C, A, G for themselves; M for A or C; R for A or G; W for A or T; S for C or G; Y for T or C; K for G or T; V for A, C, or G; H for A, C, or T; D for A, G, or T; B for C, G, or T; N for any base.

[0061] Conventionally, DNA sequences are written from 5' to 3', left-to-right.

**anti-sense DNA: 5' ATG CTT TTC ... 3'**

**sense DNA: 3' TAC GAA AAG ... 5'**

**mRNA: 5' AUG CUU UUC ... 3'**

**protein: M - L - F -**

We will use the convention that the "sense" strand is the strand used as template for mRNA synthesis.

[0062] In the present invention, the words "grow", "growth", "culture", and "amplification" mean increase in number, not increase in size of individual cells. In the present invention, the words "select" and "selection" are used in the genetic sense; i.e., a biological process whereby a phenotypic characteristic is used to enrich a population for those organisms displaying the desired phenotype.

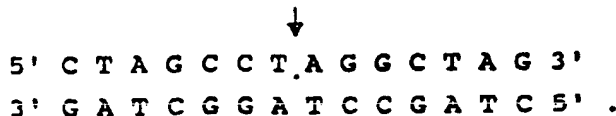
[0063] One selection is called a "selection step"; one pass of variegation followed by as many selection steps as are needed to isolate a successful DBP, is called a "variegation step". The amino acid sequence of one successful DBP from one round becomes the parental potential-DBP to the next variegation step. We perform variegation steps iteratively until the desired affinity and specificity of DNA-binding between a successful DBP and chosen target DNA sequence are achieved.

[0064] In a "forward selection" step, we select for the binding of the PDBP to a target DNA sequence; in a "reverse selection" step, for failure to bind. The target DNA sequence may be the final target sequence of interest, or the immediate target may be a related sequence of DNA (e.g., a "left symmetrized target" or "right symmetrized target"). There is an important distinction between screening and selection. Screening merely reveals which cells express or contain the desired gene. Selection allows desired cells to grow under conditions in which there is little or no growth of undesired cells (and preferably eliminates undesired cells).

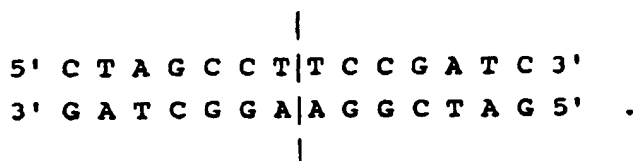
[0065] The term "operon" is used to mean a collection of one or more genes that are transcribed together. We will use operon to refer also to one or more genes that are transcribed together in eukaryotic cells independent of post-transcriptional processing.

[0066] The term "binding marker gene" is used to mean those genes engineered to detect sequence-specific DNA binding, as by association of a target DNA with a structural gene and expression control sequences. A single operon may include more than one binding marker gene (e.g., *galT<sub>K</sub>*). A "control marker gene" is one whose expression is not affected by the specific binding of a protein to the target DNA sequence. The "control promoter" is the promoter operably linked to the control marker gene.

[0067] Palindrome, palindromic, and palindromically are used to refer to DNA sequences that are the same when read along either strand, e.g.

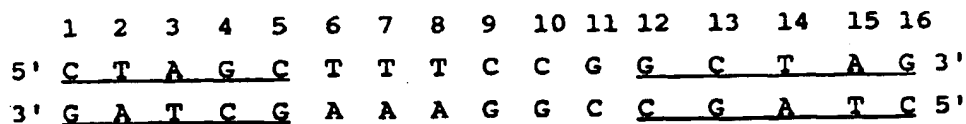
**Palindromic DNA****Rotational axis**

The arrow indicates the center of the palindrome; if the sequence is rotated 180° about the central dot, it appears unchanged. In the present application, "Palindromic" does not apply to sequences that have mirror symmetry within one strand, such as

**Mirror Plane**

DNA sequences can be partially palindromic about some point (that can be either between two base pairs or at one base pair) in which case some bases appear unchanged by a 180° rotation while other bases are changed.

[0068] A special case of partially palindromic sequence is a "gapped palindrome" in which palindromically related bases are separated by one or more bases that lack such symmetry:

**Gapped Palindrome**

has CTAGC (bases 1-5) palindromically related to GCTAG (bases 12-16) while the sequence TTTCCG (bases 6-11) in the center has no symmetry.

[0069] For the purposes of this invention, a "non-deleterious cloning site" is a region on a plasmid or phage that can be cut with one restriction enzyme or with a combination of restriction enzymes so that a large linear molecule comprising all essential elements can be recovered.

**Overview: Standard Methods**

[0070] Bacterial strains are cultured by standard methods (DAVI80, MILL72, AUSU87). Constructions of vectors are by standard methods (MANI82, ZOLL84, AUSU87). All genetic constructions are confirmed, first by analysis with restriction enzymes, and then by sequencing. Sequencing is by the Sanger dideoxy method or by Maxam Gilbert chemical method. Constructions that confer a phenotype are tested for display of the desired phenotype. These necessary controls are not described repeatedly.

**Overview: The Selection System**

[0071] The present invention separates mutated genes that specify novel proteins with desirable sequence-specific DNA-binding properties from closely related genes that specify proteins with no or undesirable DNA-binding properties.

by: 1) arranging that the product of each mutated gene be expressed in the cytoplasm of a cell carrying a chosen DNA target subsequence, and 2) using genetic selections incorporating this chosen DNA target subsequence to enrich the population of cells for those cells containing genes specifying proteins with improved binding to the chosen target DNA sequence.

[0072] A selectable deleterious gene is positioned relative to, usually downstream from, the target sequence so that the gene is not expressed if a successful DNA-binding protein specific to this target is expressed in the cell and binds the target sequence. Alternatively, a selectable beneficial gene may be arranged so that its transcription is occluded by a strong promoter (ADHY82, ELLE89a, ELLE89b). The target sequence is placed in or near the occluding promoter so that successful binding by a protein will repress the occluding promoter and allow transcription of the beneficial gene. Elledge and coworkers disclose that such systems work best in the bacterial chromosome or on low-copy-number plasmids. The cell will survive exposure to the selective conditions transcription of the selectable deleterious genetic element is blocked.

[0073] The preferred cell line or strain is easily cultured, has a short doubling time, has a large collection of well characterized selectable genes, includes variants that are deficient in genetic recombination, and has a well developed transformation system that can easily produce at least  $10^7$  independent transformants/ $\mu$ g of DNA. Bacterial cells are preferred over yeasts, fungi, plant, or animal cells because they are superior on every count. Among bacteria, *E. coli* is the premier candidate because of the wealth of knowledge of genetics and cellular processes. Other bacterial strains, such as *S. typhimurium*, *Pseudomonas aeruginosa*, *Klebsiella aerogenes*, *Bacillus subtilis*, or *Streptomyces coelicolor* could be used. DBPs that bind to host regulatory sequences, such as promoters, will be toxic. Thus, development of a DBP that specifically binds to *E. coli* promoters is preferably done in a cell line or strain, such as *S. coelicolor*, having significantly different promoter sequences.

[0074] In the most preferred embodiments, all novel DBPs are developed in *E. coli* *recA*<sup>-</sup> strains. The *recA*<sup>-</sup> genotype is preferred over other *rec*<sup>-</sup> mutations because *recA*<sup>-</sup> mutation reduces the frequency of recombination more than other known *rec*<sup>-</sup> mutations and the *recA*<sup>-</sup> mutation has fewer undesirable side effects. We choose a host strain that methylates or does not methylate the target sequence in the desired way. For example a Dcm<sup>-</sup> strain is appropriate if the target sequence contains CCwGG and we want a DBP that binds the unmethylated form.

[0075] As vectors, phage, such as M13, have the advantage of a high infectivity rate. Organisms or phage having a phase in their life cycle in which the genome is single-stranded DNA have a higher mutation rate than organisms or phage that have no phase in which the genome is single-stranded DNA. Plasmids are, however, preferred because genes on plasmids are much more easily constructed and altered than are genes in the bacterial chromosome and are more stable than genes borne on phage, such as M13. M13 derived vectors are nearly as preferred as plasmids.

[0076] In some embodiments, the cloning vector will carry: a) the selectable genes for successful DBP isolation, b) the *pdbp* gene, c) a plasmid origin of replication, and d) an antibiotic resistance gene not present in the recipient cell to allow selection for uptake of plasmid. Preferably the operative vector is of minimum size.

[0077] Alternatively, the selectable binding marker genetic elements are placed on a vector different from but compatible with the vector that carries the *pdbp* gene. This arrangement has the advantages that engineering the *pdbp* gene is easier on a smaller plasmid and manipulation of *pdbp* can not introduce mutations into the selectable binding marker genes.

[0078] Standard selections for plasmid uptake and maintenance in *E. coli* include use of antibiotics (e.g., ampicillin (Ap)) as shown in Table 2. Selection of cells with antibiotics is preferred to nutritional selections, e.g., TrpA<sup>+</sup>, for several reasons. Nutritional selection may be overcome by large volumes of cells or growth medium; host chromosomal auxotrophy is rarely total; crossfeeding of the non-growing cells by prototrophic recipients obscures the outlines of the colonies; and late mutations to prototrophy may arise on the plate due to spontaneous mutation of nongrowing cells. Nonetheless, nutritional selection may be employed.

[0079] Similarly, plasmids for use in *B. subtilis* are engineered for selection of uptake and maintenance using antibiotics. Plasmids used in streptomycete species bear genes for resistance to antibiotics such as thiostrepton, neomycin, and methylenomycin, in preference to auxotrophic markers or sporulation and pigment screens such as *spo* in bacilli and *mel* in streptomycetes.

[0080] Recombinant DNA manipulations in yeasts have been achieved using complementation of auxotrophic markers, some of which are shown in Table 3. High backgrounds are surmounted by use of two unrelated binding marker genes carried on the same vector, e.g., Leu2<sup>+</sup> and Ura3<sup>+</sup>. Selection for G418 resistance conferred by the bacterial *aphII* gene expressed in yeast offers the advantages of reduced background and a wider range of appropriate recipient strains. The current upper range of efficiency of DNA uptake into yeast cells indicates that this organism is not now preferred for the process described in this patent, although results could be achieved by large scale practice.

[0081] The selection systems must be so structured that other mechanisms for loss of gene expression are much less likely than the desired result, repression at the target DNA subsequence. Other mechanisms that could yield the desired phenotype include: point mutations that inactivate the deleterious gene or genes, deletion of the deleterious gene or genes, host mutations that suppress the deleterious genes, and repression at a site other than the target DNA

sequence.

[0082] A wide range of selectable phenotypes for *E. coli* and *S. typhimurium* have been described (VINO87). Two broad classes of selections are useful in this invention, nutritional and chemical. Such selections are inherently conditional in that they employ addition of a growth-inhibitory chemical to the selective medium, or manipulation of the nutrient components of the selective medium. Further conditionality of the preferred method is imposed by transcriptional regulation (e.g. by IPTG in combination with the *lacUV5* promoter and the *LacI<sup>q</sup>* repressor) of the variegated *pdbp* gene. In those members of the population that express DBPs that bind to the target, IPTG indirectly controls the selectable genes; in these cells, increased IPTG leads to reduced expression of the selectable genes. Therefore the phenotypes for selection are distinguished only in the presence of an inducing chemical, and potential deleterious effects of these phenotypes are avoided during storage and routine handling of the strains.

[0083] Selection of mutant strains capable of producing proteins that can bind to the target DNA subsequence is enabled by engineering conditional lethal genes or growth-inhibiting genes located downstream from the promoter that contains the target DNA subsequence. In the preferred embodiment, at least two independent conditional lethal or inhibitory selections are performed simultaneously. It is possible to use a single selection to achieve the same purpose, but this is not preferred. Two selections are strongly preferred since a simple mutation in the selected gene, occurring at a frequency of  $10^{-6}$  to  $10^{-9}$ /cell, would occur in two selected genes simultaneously at the product of the individual frequencies,  $10^{-12}$  to  $10^{-16}$ . Thus use of two selections substantially reduces the probability of isolation of artifactual revertant or suppressor strains.

[0084] Selectable genes for which both forward and reverse selections exist are preferred because, by changing host or media, we can use these genes to select for binding by a DBP to a target DNA sequence such that expression of one of these genes is repressed, or we can select phenotypes characteristic of cells in which there is no binding of the DBP. For example, expression of the *tet* gene is essential in the presence of tetracycline. On the other hand, expression of the *tet* gene is lethal in the presence of fusaric acid. Expression of the *galT* and *galK* genes in a *GalE<sup>-</sup>* host in the presence of galactose is lethal (NIKA61). Expression of *galT* and *galK* in a host that is *GalE<sup>+</sup>* and either *GalT* or *GalK<sup>-</sup>* renders the cells *Gal<sup>+</sup>* and allows them to grow on galactose as sole carbon source.

[0085] The term "source of a selective agent" includes the selective agent itself and any media components which cause the cell to manufacture the selective agent.

[0086] The Detailed Examples describe selection of strains with successful DBP binding to novel target subsequences due to turn off of two genes, each of which, if expressed, confers sensitivity to a toxic substance. It is also possible to use selection of strains in which successful DBP binding to novel target operators turns off repressors of genes encoding required gene products. For example, using the binding marker gene P22 *arc*, we place an Arc operator site so that binding of Arc represses expression of a beneficial or conditionally essential gene, such as *amp<sup>r</sup>*. Another alternative is selection of expression of required gene products due to successful binding of DBP proteins derived from positive effectors as the DBP, e.g. CAP from *E. coli*, the repressor from phage  $\lambda$ , or the Cro67 (BUSH88) mutant of  $\lambda$  Cro. Another alternative is to place the target sequence in or near a strong promoter that occludes transcription of a conditionally essential gene (ELLE89a,b).

[0087] The selections described in the Detailed Examples employ commercially available cloned genes on plasmids in strains that can be obtained from the ATCC (Rockville, MD). Alternatively, the genes can be produced synthetically from published sequences or isolated from a suitable genomic or cDNA library.

[0088] Numerous types of selections are possible for selection of DBP expression in *E. coli*. The toxic and inhibitory agents listed in Table 4 are used with appropriately engineered host strains and vectors to select loss of gene function listed above. Repression of transcription of these genes allows growth in the presence of the agents. Other outcomes such as deletions or point mutations in these genes may also be selected with these agents, hence two functionally unrelated selections are used in combination. These agents share the property that cell metabolism is stopped, and unlike the nutritional selections, the inhibitory agents are not overcome by components of the growth medium or turnover of macromolecules in the cells. Selections using antibiotics, metabolite analogs, or inhibitors are preferred. Another class of selections includes those for repression of phage or colicin receptors, or for repression of phage promoters. These agents kill by single-hit kinetics, and in the case of phage, are self-replicating, making the multiplicity of agent to putative repressed cell much more difficult to control and so are not preferred (BENS86).

[0089] Any selection system relevant to the cell line or strain may be substituted for those in the examples given here, with appropriate changes in the engineering of the cloning vectors. One example is the dominant *pheS<sup>+</sup>* gene carried on plasmid pHE3 (ATCC #37.161) in a *pheS12* background. Turn-off of *pheS<sup>+</sup>* is selected with p-fluorophenylalanine (Sigma Corp., St. Louis, MO).

[0090] We could choose the *Streptomyces coelicolor* cloned glucose kinase gene for selection of the DBP<sup>+</sup> phenotype, using the metabolite analog deoxyglucose.

[0091] Each batch of antibiotic is checked for MIC (minimum inhibitory concentration) under the condition of use. Increased concentration of antibiotic may be used to increase the stringency of the selection, in most cases.

[0092] The user varies the medium formulation (pH, cation concentrations, buffering agent, etc.) for a particular

selection if the results are not optimal with the strain at hand. For example, Maloy and Nunn (MALO81) describe a medium yielding improved selection of *Fus<sup>R</sup> E. coli* colonies from a *Tc<sup>R</sup>* background, compared to the medium employed by Bochner (BOCH80) for this purpose using *S. typhimurium*.

[0093] Stringency of selection can be modulated by controlling copy number of plasmids bearing the selectable genes; increasing copy number of selectable genes increases the stringency of the selection.

[0094] During the initial phases of the progressive development of DBP molecules, it is desirable to produce a high intracellular concentration of DBP. The stringency of the selection is increased in subsequent phases of successful DBP development by allowing fewer molecules of DBP per cell. Thus it is preferred to regulate transcription of *pdbp* by an inducible or derepressible promoter, such as *PlacUV5*.

[0095] High total cell input often decreases stringency of selections, by providing metabolites that are specifically omitted, by mass action with respect to an inhibitory agent, or by generating a large number of artificial satellite colonies that follow the appearance of genetically resistant colonies. The number of cells that are successfully transformed is a function of efficiency of ligation and transformation processes, both of which are optimized in the embodiment of this invention. Procedures for maximal transformation and ligation efficiency are from Hanahan (HANA85) and Legerski and Robberson (LEGE85) respectively. Increasing stringency is imposed under the conditions of high efficiency of these processes by inoculation of plates with small volumes or dilutions of cell samples. Pilot experiments are performed to determine optimum dilution and volume.

[0096] In Detailed Example 1, the transformation event is followed by dilution and growth of cells in permissive medium following transformation. Exogenous inducer of DBP expression is included at this step, and a set of selections are then imposed in liquid medium. Surviving cells are concentrated by centrifugation, and selected for these and additional traits using solid medium in Petri plates. This protocol offers the advantage that fewer identical siblings are obtained and a larger population is easily screened. In Detailed Example 1, repression of the *Gal<sup>S</sup>* phenotype is selected by exposing transformants to galactose in liquid medium, which produces visible lysis of galactose sensitive cells. The second selection employed in Detailed Example 1 is for the *Fus<sup>R</sup>* phenotype due to repression of *Tc<sup>R</sup>*, which requires limitation of total inoculum size to  $10^6$  cells/plate. Similar protocol variations are introduced to combine selections for transformation and successful DBP function.

[0097] Tests of selective agents to determine the conditions that kill or inhibit sensitive cells are performed with pure cultures of sensitive cells. These include strains carrying the selective marker genes having the recognition sequence of the IDBP as target, with and without *idbp*, and with and without the inducer of *idbp* expression.

[0098] Cultures of sensitive cells are applied to selective media as inocula appropriate to the selection (usually  $10^6$  to  $10^8$  per plate). Sufficient numbers of replicates ( $10^7$  to  $10^9$  total sensitive cells for each medium) are tested by each selection. The rate at which the cultures produce revertants and phenotypic suppressors (considered together as revertants) is determined. A rate greater than  $10^{-6}$  per cell indicates that stringency must be increased. If reversion rates are below this level, as we have shown for the selections described in Example 1, mixing experiments are performed to determine the sensitivity of recovery of a small fraction of resistant cells from a vast excess of sensitive cells.

[0099] Normally, the deleterious gene product of a binding marker gene is a protein. It may also be an RNA, e.g., an mRNA which is antisense to the mRNA of an essential gene and therefore blocks translation of the latter mRNA into protein. Another alternative is that transcription of the binding marker gene may be deleterious because this transcription occludes transcription of an adjacent beneficial gene. Selectively deleterious genes suitable for use in the present invention include those shown in Table 4.

[0100] The two selectably deleterious genes are preferably not functionally related. For example, the chosen genes should not code for proteins localized to or affecting the same macromolecular assembly in the cell or which alter the same or intersecting anabolic or catabolic pathways. Thus, use of two inhibitors that select for mutations affecting RNA synthesis, aromatic amino acid synthesis, or each of histidine and purine synthesis are not preferred. Similarly, two inhibitors that are transported into the cell by shared membrane components are thus functionally related, and are not preferred. In this manner the user reduces the frequency of isolation of single host mutations that yield the apparent desired phenotype, because of suppression of the shared functionality, interacting component, or precursor relationship. Host mutations of this type are conveniently distinguished by a screen of the selectable phenotypes in the absence of the inducer of the DBP, e.g. IPTG.

[0101] Examples of pairs of deleterious genes which are recommended for use in the present invention are given in Table 5A. In each case, one of the paired genes codes for a product that acts intracellularly while the other codes for a product that acts either in transport into or out of the cell or acts in an unrelated biological pathway. Table 5B gives some pairs that are not recommended. These pairs have not been shown to malfunction, but they are not recommended, given the large number of choices that are clearly functionally unrelated.

[0102] A preferred novel feature is the use of a copy of the promoter of one of these beneficial or conditionally essential genes, operably linked to the target DNA subsequence, to direct transcription of the selectably deleterious or conditionally lethal binding marker genes of the plasmid. If the potential-DBP should repress the selectable gene by binding to this promoter, it would also repress this beneficial activity.

[0103] In order to assure that selection for DBP binding is specific to the target and not the promoter, we, preferably, place one of the two selectable binding marker genes under the same transcription initiation signal as the gene we use for selection of vector maintenance. In Detailed Example 1, transcription of the *galT* and *galK* genes is initiated by the  $P_{amp}$  promoter, as is the *amp* gene.

5 [0104] It is possible that the potential-DBP will bind specifically to the boundary between the target DNA sequence and the promoter, or within the structural gene. In the preferred embodiment, we discriminate against this mechanism by choosing a different promoter, operably linked to another copy of the same target DNA sequence, for the second selectable gene. Preferably, the two promoters that initiate transcription of the selectable genes should be strong enough to give a sensitive selection, but not too strong to be repressed by binding of a novel DBP. Some well studied  
10 promoters and their scores by the Mulligan algorithm (MULL84) are shown in Table 6. Promoters that score between 50% and 70% are good candidates for use in binding marker genes. Preferably, the two promoters have significant sequence differences, particularly in the region of the junction to the target DNA sequence. Specifically, the region between the -10 region and the target sequence, which comprises five to seven bases, should have no more than two identical bases in the two promoters. Although the -10 regions of promoters show high homology, promoters are known  
15 (e.g.  $P_{amp}$  having GACAAT and  $P_{neo}$  having TAAGGT) that have as few as two out of six bases identical in this region, and such difference is preferred.

[0105] The target DNA sequence for the potential DNA-binding protein must be associated with the two deleterious binding marker genes and their promoters so that expression of the binding marker genes is blocked if a novel protein in fact binds to the target sequence. The target DNA sequence could appear upstream of the gene, downstream of the  
20 gene, or, in certain hosts, in a noncoding region (*viz.* an "intron") within the gene. Preferably, it is placed upstream of the coding region of the gene, that is, in or near the RNA polymerase binding site for the gene, *i.e.* the promoter. If the binding marker gene is an occluding promoter, the target is, preferably, placed downstream of the promoter. Placement of the target DNA sequence relative to the promoter is influenced by two main considerations: a) protein binding should have a strong effect on transcription so that the selection is sensitive, b) the activity of the promoter in the absence of a  
25 binding protein should be relatively unaffected by the presence of the test DNA sequence compared to any other target subsequence.

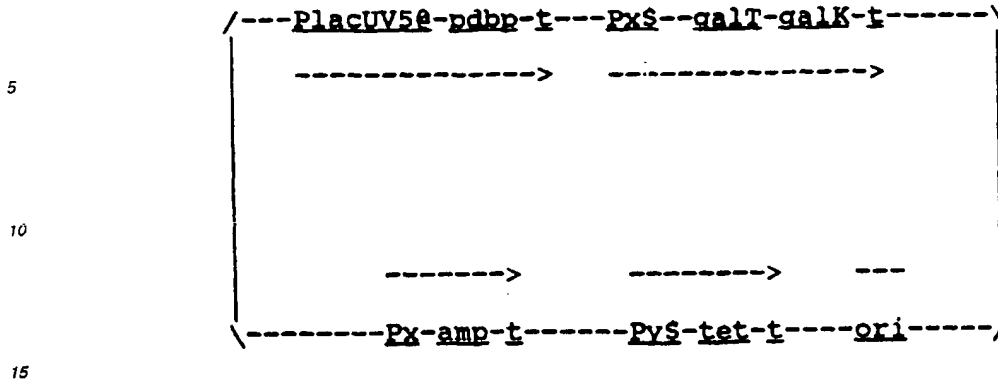
[0106] In the present invention, we will deal primarily with DNA target subsequences of 10 to 25 bases. It has been noted that the highly conserved -35 region and the highly conserved -10 region are separated by between 15 and 21 base pairs with a mode of 17 base pairs (HAWL83, MULL84). Some of the bases between -35 and -10 are statistically  
30 non-random; thus placement of target DNA sequences longer than 10 bases between the -10 and -35 regions would likely affect the promoter activity independent of binding by potential-DBPs. Because quantitative relationships between promoter sequence and promoter strength are not well understood; it is preferable, at present, to use known promoters and to position the target at the edge of the RNA polymerase binding site.

[0107] Protein binding to DNA has maximum effect on transcription if the binding site is in or just down-stream from the promoter of a gene. Hoopes and McClure (HOOP87) have reviewed the regulation of transcription initiation and report that the LexA binding site can produce effective repression in a variety of locations in the promoter region. In a preferred embodiment, we place the target DNA sequences that begin with A or G so that the first 5' base of the target sequence is the +1 base of the mRNA, as the LexA binding site is located in the *uvrD* gene (HOOP87, p1235). If the target sequence begins with C or T, we preferably place the target so that the first 5' base of the target is the +2 base  
40 of the mRNA and we place an A or G at the +1 position. An alternative is to place the target DNA sequences upstream of the -35 region as the LexA binding site is located in the *ssb* gene (HOOP87, p1235).

[0108] It may be useful in early stages of the development of a DBP to have more than one copy of the target DNA sequence positioned so that binding of a DBP reduces transcription of the selectable gene. Multiple copies of the target DNA sequence enhances the sensitivity of phenotypic characteristics to binding of DBPs to the target DNA sequence.  
45 Multiple copies of the target DNA sequence are, preferably, placed in tandem downstream of the promoter. Alternatively, one could place one copy upstream of the promoter and one or more copies downstream.

[0109] We arrange the genes on the plasmid or plasmids in such a way that no single deletion event eliminates both deleterious genes without also eliminating a gene essential either to plasmid replication or cell survival. Thus, resistant colonies are unlikely to arise through deletions because two independent deletion events are required. Similarly, simultaneous occurrence of one point mutation and one deletion is as unlikely as two point mutations or two deletions.  
50

[0110] A typical arrangement of genes on the operative cloning vector, similar to that used in Detailed Example 1, is:



P<sub>x</sub> represents the promoter that initiates transcription of the amp gene. A second copy of P<sub>x</sub> initiates transcription of galT,K. P<sub>y</sub> is a promoter driving tet, t is a transcriptional terminator (different terminators may be used for different genes), and \$ is the target subsequence. PlacUV5 is the lacUV5 promoter, @ represents the lacO operator, and pdpb is a variegated gene encoding potential DBPs. Placement of the pdpb relative to other genes is not important because mutations or deletions in pdpb cannot cause false positive colony isolates. Indeed, it is not necessary that the pdpb gene be on the selection vector at all. The purpose of the selection vector is to ensure that the host cell survives only if the one of the PDBPs binds to the target sequence (forward selection) or fails to so bind (reverse selection). The pdpb gene may be introduced into the host cell by another vector.

[0111] Two-way selections are available for both tet and galT,K (vide supra). The orientation of each gene in the selection vector is unimportant because strong terminators (e.g. rrnBt1, rrnBt2, phage fd terminator) are preferably placed at the ends of each transcription unit. That galT,K and tet are separated by essential genes, however, is of fundamental importance. The sequence ori is essential for plasmid replication, and the amp gene, the transcription of which is initiated by P<sub>x</sub>, is essential in the presence of Ap. Successful repression of galT,K and tet is selected with galactose and fusaric acid. No single deletion event can remove both the latter genes and allow plasmid maintenance or cell survival under selection. In addition, binding by a novel DBP to the P<sub>x</sub> promoter would render the cell Ap sensitive. These arrangements make appearance of a novel DBP that binds the target DNA more probable than any of the other modes by which the cells can escape the designed selections.

#### Overview: Choice of target DNA binding sequence for development of successful novel DBPs:

[0112] Our goal is the development, in part by conscious design and in part by in vivo selection, of a protein which binds to a DNA sequence of significance, e.g., a structural gene or a regulatory element, and through such binding inhibits or enhances its biological activity. In the preferred embodiment, the protein represses transcription of a deleterious element, such as a viral gene. A sufficiently long sequence could be the target of several independently acting DBPs.

[0113] Another goal of this invention is to derive one or more DBPs that bind sequence-specifically to any predetermined target DNA subsequence. It is not yet possible to design the DBP-domain amino-acid sequence from a set of rules appropriate to the target DNA subsequence. Rather, it is possible to pick sets of residues that can affect the DNA recognition of a parental DBP. Then, variegation of residues that affect DNA recognition coupled with selection for binding to the target DNA subsequence can produce a novel DBP specific for the target DNA subsequence. Such a method is limited by the number of amino acids that can be varied at one time. To develop a novel DBP that recognizes 15 bases could require changing 15 or more residues in the initial DBP. Variegation of 15 residues through all 20 amino acids would produce  $20^{15} = 3.3 \times 10^{19}$  sequences and is beyond current technology. Thus we start with the recognition sequence of the initial DBP, change two to five bases and select, in one or more rounds of variegation and selection, a novel DBP that recognizes this new target DNA subsequence. This new DBP becomes the parent to the next step in which the target DNA subsequence is changed by an additional two to five bases so that a stepwise series of changes in binding protein and changes in target is used. It is emphasized here that, although we initially select DBPs that recognize sequences similar to that recognized by the IDBP, the ultimate target sequence recognized by the desired final DBP can be completely unrelated to the recognition sequence of the IDBP.

[0114] The process of finding a DBP that recognizes a sequence within a genome is shortened if we pick sequences that have some similarity to the cognate sequence of the initial DBP. The intent is to locate several unique sites in the gene which can be bound specifically by DBPs such that transcription through those sites is reduced.



[0115] The sequences of some regions of genes of eukaryotic pathogens vary among strains (SAAG88). To optimize the search for target sites in the gene selected for repression such that repression will be effective in all or the majority of strains of a pathogen, regions of conserved DNA sequence within the gene are, preferably, identified.

[0116] There may be a very small number of sequences that occur in the genome of the host cells for which binding of a DBP will be lethal. For this reason, the regulatory sequences, such as promoters, of the host organism are not preferred targets for DBP development. Preferably, the target sequence occurs only in the gene of interest. For some applications, target sequences that occur at locations other than the site of intended action may be used if binding of a protein to the extra sites is acceptable.

[0117] Preliminary elimination of non-unique sequences is done by searching DNA sequence data banks of host genomic sequences and bacterial strain sequences; and by searching the plasmid sequences for matches to the potential target subsequences. Remaining potential target subsequences are then used as oligonucleotide probes in Southern analyses of host genomic DNA and bacterial DNA. Sequences which do not anneal to host or bacterial DNA under stringent conditions are retained as target subsequences. These target subsequences are cloned into the operative vector at the promoters of the selection genes for DBP function, as described for the test DNA binding sequence.

[0118] Choice of target subsequences is based also on the optimal location of target sites within a gene such that transcription will be maximally affected. Studies of monkey L-cells show that *lac* repressor can bind to *lac* operator, or to two *lac* operators in tandem, in the L-cell nucleus (HUMC87, HUMC88). Further, this binding results in repression of a downstream chloramphenicol acetyl transferase gene in this system, and repression is relieved by IPTG. Two tandem operators repress CAT enzyme production to a greater extent than a single operator. The user preferably locates two to four target sites relatively close to each other within the transcriptional unit.

#### Overview: Selection of the Initial DNA-Binding Protein for Variegation

[0119] The choice of an initial DBP is determined by the degree of specificity required in the intended use of the successful DBP and by the availability of known DBPs. The present invention describes three broad alternatives for producing DBPs having high specificity and tight binding to target DNA sequences. The present invention is not limited to these classes of initial potential DBPs.

[0120] A first alternative is to use a polypeptide that will conform to the DNA and can wind around the DNA and contact the edges of the base pairs. A second alternative is to use a globular protein (such as a dimeric H-T-H protein) that can contact one face of DNA in one or more places to achieve the desired affinity and specificity. A third alternative is to use a series of flexibly linked small globular domains that can make contact with several successive patches on the DNA.

#### DNA features influencing choice of an initial DBP:

[0121] Features of DNA that influence the choice of an initial DBP include sequence-specific DNA structure and the size of the genome within which the DBP is expected to recognize and affect gene expression.

[0122] Sequence-specific aspects of DNA structure that can influence protein binding include: a) the edges of the bases exposed in the major groove, b) the edges of the bases exposed in the minor groove, c) the equilibrium positions of the phosphate and deoxyribose groups, d) the flexibility of the DNA toward deformation, and e) the ability of the DNA to accept intercalated molecules. Note that the sequence-specific aspects of DNA are carried mostly inside a highly charged molecular framework that is nearly independent of sequence.

[0123] The strongest signals of sequence are found in the edges of the base pairs in the major groove, followed by the edges in the minor groove. The groove dimensions depend on local DNA sequence (NEID87b, KOUD87, ULAN87).

[0124] The number of base pairs required to define a unique site depends on the size and non-randomness of the genome. Consider a genome of length  $Z_g$  bases and consider a specific subsequence of length  $Q$ . If the genome is random, the subsequence is expected to occur  $N(Q)$  times, where

$$N(Q) = \frac{2^{Z_g}}{4^Q} \approx \frac{2^{Z_g}}{2^{2Q}}$$

From this equation, we derive the expression  $Q_{\min}$ , which is the lower limit of the length of subsequences that are expected to occur once or be absent:

$$Q_u = \log_2(2 Z_g)/2.$$

| $Z_g$     | $\log_2(2 Z_g)/2$ | $Q_u$ |
|-----------|-------------------|-------|
| $10^6$    | 10.5              | 11    |
| $10^7$    | 12.1              | 13    |
| $10^8$    | 13.8              | 14    |
| $10^9$    | 15.5              | 16    |
| $10^{10}$ | 17.1              | 18    |

[0125] Thus, a DNA subsequence comprising 12 base pairs may be unique in the *E. coli* genome ( $5 \times 10^6$  bp), but is likely to occur about 180 times in a random sequence the size of the human genome ( $3 \times 10^9$  bp).

[0126] The non-random nature of DNA sequences in genomes has been shown to result in the over- and under-representation of specific sequences. The random-genome model can under-estimate the probe length needed to define a unique coding sequence (LATH85). Recognition sites for certain restriction enzymes occur in clusters and are found much more often than expected (SMIT87). In contrast, *lac* repressor binding sites in eukaryotic genomes are almost two orders of magnitude less frequent than expected on the basis of random sequence (SIM084).

#### Protein features influencing choice of initial DBP:

[0127] Sequence-specific binding to DNA by DBPs does not require unpairing of the bases. Most sequence-specific binding by proteins to DNA is thought to involve contacts in the DNA major groove.

[0128] To be certain of unique recognition in the human genome, it is best to design a protein that recognizes 19 to 21 base pairs. To contact 20 base pairs directly, a protein would need to: a) wind two full turns around the DNA making major groove contacts, b) make a combination of major groove and minor groove contacts, or c) contact the major groove at four or five places. An extended polypeptide, binding in the major groove of B-DNA, lies about 5.0 Å from the DNA axis. One base pair and 1 1/2 amino acids extend roughly equal distances along the helix (SAEN83, p238).

[0129] A nine residue alpha helix, such as the recognition helices of H-T-H repressors, extends about 13.5 Å along the major groove. If residues with long side chains are located at each terminus of the helix, the helix can make contacts over a 20.0 Å stretch of the major groove allowing six base pairs to be contacted. Parts of the DBP other than the second helix of the H-T-H motif can make additional protein-DNA contacts, adding to specificity and affinity. The rigidity of the alpha helix prevents a long helix from following the major groove around the DNA. A series of small domains, appropriately linked, could wind around DNA, as has been suggested for the zinc-finger proteins (BERG88a, GIBS88, FRAN88). In an extended configuration a polypeptide chain progresses roughly 3.2 to 3.5 Å between consecutive residues. Thus, a 10 residue extended protein structure could contact 5 to 8 bases of DNA.

[0130] Stable complexes of proteins with other macromolecules involve burial of  $1000 \text{ Å}^2$  to  $3000 \text{ Å}^2$  of surface area on each molecule. For a globular protein to make a stable complex with DNA, the protein must have substantial surface that is already complementary to the DNA surface or can be deformed to fit the surface without loss of much free energy. Considering these modalities we assign each genetically encoded polypeptide to one of three classes:

1) a polypeptide that can easily deform to complement the shape of DNA,

2) a globular protein, the internal structure of which supports recognition elements to create a surface complementary to a particular DNA subsequence, and

3) a sequential chain of globular domains, each domain being more or less rigid and complementary to a portion of the surface of a DNA subsequence and the domains being linked by amino acid subsequences that allow the domains to wind around the DNA.

[0131] Complementary charges can accelerate association of molecules, but they usually do not provide much of the free energy of binding. Major components of binding energy arise from highly complementary surfaces and the liberation of ordered water on the macromolecular surfaces.

Properties of sequence-specific DNA-binding by polypeptides:

[0132] An extended polypeptide of 24 amino acids lying in the major groove of B-DNA could make sequence-specific interactions with as many as 15 base pairs, which is about the least recognition that would be useful in eukaryotic systems. Peptides longer than 24 amino acids can contact more base pairs and thus provide greater specificity.

[0133] Extended polypeptide segments of proteins bind to DNA in natural systems (e.g.  $\lambda$  repressor and Cro, P22 Arc and Mnt repressors). The DNA major groove can accommodate polypeptides in either helical or extended conformation. Side groups of polypeptides that lie in the major groove can make sequence-specific or sequence-independent contacts. Since the polypeptide can lie entirely within the major groove, contacts with the phosphates are allowed but not mandatory. Thus a polypeptide need not be highly positively charged. A neutral or slightly positively charged polypeptide might have very low non-specific binding.

[0134] Polypeptides composed of the 20 standard amino acids are not flat enough to lie in the minor groove unless the sequence contains an extraordinary number of glycines, however, residue side-groups could extend into the minor groove to make sequence-specific contacts. Polypeptides of more than 50 amino acids may fold into stable 3D structures. Unless part of the surface of the structure is complementary to the surface of the target DNA subsequence, formation of the 3D structure competes with DNA binding. Thus polypeptides generated for selection of specific binding are preferably 25 to 50 amino acids in length.

[0135] Polypeptides present the following potential advantages:

- a) low molecular weight: an extended polypeptide offers the maximum recognition per amino acid,
- b) polypeptides have no inherent dyad symmetry and so are not biased toward recognition of palindromic sequences,
- c) polypeptides may have greater specificity than globular proteins, and
- d) peptides may be good models from which other low molecular weight compounds may be designed.

[0136] Thus, one would choose a polypeptide as initial DNA-binding molecule if high specificity and low molecular weight are desired.

[0137] No sequence-specific DNA-binding by small polypeptides has been reported to date. Possible reasons that such polypeptides have not been found include: a) no one has sought them, b) cells degrade polypeptides that are free in the cytoplasm, and c) they are too flexible and are not specific enough.

[0138] In a preferred embodiment, a DNA-binding polypeptide is associated with a custodial domain to protect it from degradation, as discussed more fully in Examples 3 and 4.

Properties of globular proteins influencing choice of initial DBP:

[0139] The majority of the well-characterized DBPs are small globular proteins containing one or more DNA-binding domains. No single-domain globular protein comprising 200 or fewer amino acids is likely to fold into a stable structure that follows either groove of DNA continuously for 10 bases. The structure of a small globular protein can be arranged to hold more than one set of recognition elements in appropriate positions to contact several sites along the DNA thereby achieving high specificity, however, the bases contacted are not necessarily sequential on the DNA. For example, each monomer of  $\lambda$  repressor contains two sequence-specific DNA recognition regions: the recognition helix of the H-T-H region contacts the front face of the DNA binding site and the N-terminal arm contacts the back face. To obtain tight binding, a globular protein must contact not only the base-pair edges, but also the DNA backbone making sequence-independent contacts. These sequence-independent contacts give rise to a certain sequence-independent affinity of the protein for DNA. The bases that intervene between segments that are directly contacted influence the position and flexibility of the contacted bases. If the DNA-protein complex involves twisting or bending the DNA (e.g. 434 repressor-DNA complex), non-contacted bases can influence binding through their effects on the rigidity of the target DNA sequence.

[0140] The phage repressors Arc, Mnt,  $\lambda$  repressor and Cro are proposed to bind to DNA at least partly via binding of extended segments of polypeptide chain. The N-terminal arm of  $\lambda$  repressor makes sequence-specific contacts with bases in the major groove on the back side of the binding site. The C-terminal "tail" of  $\lambda$  Cro is proposed to make sequence-independent contacts in the minor groove of the DNA. The structure of neither Arc nor Mnt has been determined; however, the sequence specificity of the N-terminal arm of Arc can be transferred to Mnt; viz. when Arc residues 1-9 are fused to Mnt residues 7 through the C-terminal, the fusion protein recognized the arc operator but not the mnt operator. Residues 2, 3, 4, 5, 8, and 10 of Arc have been proposed to contact operator DNA and residue 6 of Mnt has

been shown to be involved in sequence-specific operator contacts.

[0141] Binding to non-palindromic sequences requires alteration of dyad-symmetric proteins. Even non-palindromic DNA has approximate dyad symmetry in the deoxyribosephosphate backbone; proteins that are heterodimers or pseudo-dimers engineered from known globular DBPs are good candidates for the mutation process described here to obtain globular proteins that bind non-palindromic DNA. It has been observed that the DNA restriction enzymes having palindromic recognition are composed of dyad symmetric multimers (MCCL86), while restriction enzymes and other DNA-modifying enzymes (e.g. Xis of phage  $\lambda$ ) having asymmetric recognition are comprised of a single polypeptide chain or an asymmetric aggregate (RICH88). Such proteins may also provide reasonable starting points to generate DBPs recognizing non-palindromic sequences.

[0142] A globular protein can bind sequence-specifically to DNA through one set of residues and activate transcription from an adjacent gene through a different set of residues (for example,  $\lambda$  or P22 repressors). The internal structure of the protein establishes the appropriate geometric relationship between these two sets of residues. Globular proteins may also bind particular small molecules, effectors, in such a way that the affinity of the protein for its specific DNA recognition subsequence is a function of the concentration of the particular small molecules (e.g. CRP and cAMP). Conditional DNA-binding and gene activation are most easily obtained by engineering changes into known globular DBPs.

[0143] Some DBPs from bacteria and bacteriophage have been shown to have sufficient specificity to operate in mammalian cells.

[0144] An initial DBP may be chosen from natural globular DBPs of any cell type. The natural DBP is preferably small so that genetic engineering is facile. Preferably, the 3D structure of the natural DBP is known; this can be determined from X-ray diffraction, NMR, genetic and biochemical studies. Preferably, the residues in the natural DBP that contact DNA are known. Preferably the residues that are involved in multimer contacts are known. Preferably the natural operator of the natural DBP is known. More preferably, mutants of the natural operator are known and the effects of these mutants on binding by natural DBP and mutant DBPs are known. Preferably, mutations of the DBP are known and the effects on protein folding, multimer formation, and *in vivo* half life-time are known. Most of the above data are available for  $\lambda$  Cro,  $\lambda$  repressor and fragments of  $\lambda$  repressor, 434 repressor and Cro proteins, *E. coli* CRP and *trp* repressor, P22 Arc, and P22 Mnt.

[0145] Globular DBPs are the best understood DBPs. In many cases, globular DBPs are capable of sufficient specificity and affinity for the target DNA sequence. Thus globular DBPs are the most preferred candidates for initial DBP. Table 8 contains a list of some preferred globular DBPs for use as initial DBPs.

[0146]  $\lambda$  repressor and phage 434 repressor have been extensively studied (CHAD71, PTAS80, PABO79, JOHN79, SAUE79, SAUE86, PABO82a,b, LEWI83, OHLE83, WEIS87a,b,c, REID88, ANDE87, NELS86, ELIA85). Both proteins comprise an amino-terminal DNA-binding domain having four homologous alpha helices. Helices 2 and 3 form the H-T-H motif. DNA contacts originate in helix 2, helix 3, and adjacent regions with helix 3 providing most of the contacts. The N-terminal domains of  $\lambda$  repressor contact each other along helix 5 (PABO82b) while in 434 repressor the interdomain contacts are beyond helix 4, there being no helix 5 (ANDE87).

[0147] The operator DNA bends symmetrically in the 434 repressor-consensus operator co-crystal (ANDE87). The center of the 14 base pair DNA helix is over-wound and bends slightly along its axis such that it curls around the alpha 3 helix of each repressor monomer; the ends of the operator DNA helix are underwound. Bending of operator DNA has also been proposed in models of Cro protein and CAP protein operator binding (OHLE83, GART88). Consistent with the results of Gartenberg and Crothers, bending of the 434 operator toward Cro is toward the minor groove and occurs most readily when the central bases consist exclusively of A and T (KOU87); in this case, substitution of CG base pairs greatly reduces binding.

[0148]  $\lambda$  Cro (TAKE77) has been described from an X-ray structure of the protein without DNA (ANDE81). Alpha helix 2 lies across the operator major groove and may make contacts to operator backbone phosphates at its N-terminal and C-terminal ends. In addition, backbone phosphates may be contacted by residues at the C terminus of alpha 3, N terminus of beta 2, and C terminus of beta 3 (PABO84). In computer model building of  $\lambda$  Cro-operator DNA interactions, bending of operator DNA or bending at the monomer-monomer interface of the Cro dimer have been proposed to make the best fit between operator and dimer (PABO84).

[0149] Key amino acids within the H-T-H region of 434 Cro and  $\lambda$  Cro are highly conserved (PABO84), and 434 Cro binds operator DNA as a dimer (WHAR85a). Because the crystals of 434 Cro and DNA do not diffract to high resolution, atomic details of the protein-DNA interactions are not revealed (WOLB88). Nevertheless, Wolberger *et al.* report very significant similarities and differences between the DNA binding patterns of 434 repressor and 434 Cro. These observations on DBPs from 434, together with recent results on Trp repressor (OTWI88), support the view that a) structural elements that fit into the major groove of DNA can function in a variety of closely related ways, b) bending of DNA complexed to proteins is an important determinant of specificity, and c) that mechanisms of recognition may be quite subtle.

[0150] Crystal structures have been determined for two DBPs, CRP (WEBE87a) and TrpR (OTWI88) from *E. coli*. Both these proteins contain H-T-H motifs and bind their cognate operators only when particular effector molecules are bound to the protein, cAMP for CRP and L-tryptophan for TrpR. Binding of each effector molecule causes a conformational

tional change in the protein that brings the DNA-recognizing elements into correct orientation for strong, sequence-specific binding to DNA (JOHN86). The DNA-binding function of Lac repressor is also modulated through protein binding of an effector molecule (e.g. lactose); unlike CRP and TrpR, Lac repressor binds DNA only in the absence of the effector. CRP can act either as an activator (RENY88) or as a repressor (POLA88) depending on the relationship between the CRP-binding site and the rest of the promoter.

[0151] Two structures of CRP (MCKA81, MCKA82) and one structure of a CRP mutant (WEBE87a) are available. Otwinowski *et al.* (OTWI88) have published an X-ray crystal structure of TrpR bound to the Trp operator. This structure shows that, although TrpR contains a canonical H-T-H motif, the positioning of the recognition helix with respect to the DNA is quite different from the positioning of the corresponding helix in other H-T-H DBPs (MATT88) for which structures of protein-DNA complexes are available. Unlike previously determined structures, most of the interactions between atoms of TrpR and bases are mediated by localized water molecules. It is not possible to distinguish between localized water and atomic ions, such as Na<sup>+</sup>, by X-ray diffraction alone. We shall follow Otwinowski *et al.* and refer to these peaks in electron density as water, although ions cannot be ruled out.

[0152] Bass *et al.* (BASS88) studied the binding of wild type TrpR and single amino acid missense mutants of TrpR to a consensus palindromic Trp operator and to palindromic operators that differ from the consensus by a symmetric substitution at one base in each half operator. Bass *et al.* conclude that the contact between the H-T-H motif of TrpR and the operators must be substantially different from the model that had been built based on the 434 Cro-DNA structure.

[0153] Thus the binding of globular DBPs that are modulated by effector molecules is fundamentally the same as the binding of unmodulated globular DBPs, but the details of each protein's interactions with DNA are quite different. Prediction of which amino acids will produce strong specific binding is beyond the capabilities of current theory. Given the important role of localized waters or ions in the TrpR-DNA interface (OTWI88) and in the 434R-DNA interface (AGGA88), such predictions are likely to remain beyond reach for some time.

[0154] The Mnt repressor of P22 is an 82 residue protein that binds as a tetramer to an approximately palindromic 17 base pair operator presumably in a manner that is two-fold rotationally symmetric. Although the Mnt protein is 40% alpha helical and has some homology to  $\lambda$  Cro protein, Mnt is known to contact operator DNA by N-terminal residues (VERS87a) and possibly by a residue (K79) close to the C terminus (KNIG88). It is unlikely, therefore, that an H-T-H structure in Mnt mediates DNA binding (VERS87a). Another residue (Y78) close to the C-terminal end has been found to stabilize tetramer formation (KNIG88). Though the three dimensional structure of Mnt is not known, DNA-binding experiments have indicated that the Mnt operator, in B-form conformation, is contacted at major groove nucleotides on both front and back sides of the operator helix (VERS87a).

[0155] The Arc repressor of P22 is a 53 residue protein that binds as a dimer to a partially palindromic 21 base pair operator adjacent to the mnt operator in P22 and protects a region of the operator that is only partially symmetric relative to the symmetric sequences in the operator (VERS87b). Arc is 40% homologous to the N-terminal portion of Mnt, and the N-terminal residues of the Arc protein contact operator DNA such that an H-T-H binding motif is unlikely, as in Mnt binding (VERS86b). The three dimensional structure of Arc, like Mnt, is not known, but a crystallographic study is in progress (JORD85). DNA-binding experiments have shown that Arc probably binds along one face of B-form operator DNA. These experiments indicate that Arc contacts operator phosphates farther out from the center of operator symmetry than do the repressors or Cro proteins of  $\lambda$  or 434, or P22 Mnt protein. Thus the researchers state that the operator DNA may be bent around Arc in binding or Arc dimer may have an extended structure to allow such contacts to occur (VERS87b). These alternatives are not mutually exclusive.

#### DNA-Binding Proteins Other Than Repressor Proteins

[0156] Any protein (or polypeptide) which binds DNA may be used as an initial DNA-binding protein; the present method is not limited to repressor proteins, but rather includes other regulatory proteins as well as DNA-binding enzymes such as polymerases and nucleases.

[0157] Derivatives of restriction enzymes may be used as initial DBPs. All known restriction enzymes recognize eight or fewer base pairs and cut genomic DNA at many places. Expression of a functional restriction enzyme at high levels is lethal unless the corresponding sequence-specific DNA-modifying enzyme is also expressed. EcoRI that lacks residues 1-29, denoted EcoRI-delN29, has no nuclease activity (JENJ86); EcoRI-delN29 binds sequence-specifically to DNA that includes the EcoRI recognition sequence, GAATTC, (BECK88).

[0158] From the structure of R.EcoRI (MCCL86), we can see that extension of the polypeptide chain at either the amino or carboxy terminus would allow contacts with base pairs outside of the canonical hexanucleotide.

[0159] Specifically, extending EcoRI(AT139), EcoRI(GS140), or EcoRI(RQ203) (YANO87) by, for example, ten highly variegated residues at the amino terminus and selecting for binding to a target such as, TGAATTCA or GGAAT-TCC, allows isolation of a protein having novel DNA-recognition properties. Alternatively, EcoRI may be extended at the amino terminus by addition of a zinc-finger domain. It may be useful to have two or more tandem repeats of the octa-

nucleotide target placed in or near the promoter region of the selectable gene. Fox (FOXK88) has used DNase-I to footprint EcoRI bound to DNA and reports that 15 bp are protected. Thus, repeated octanucleotide targets for proteins derived from EcoRI should be separated by eight or more base pairs; one could place one copy of the target upstream of the -35 region and one copy downstream of the -10 region. There are many residues in EcoRI that contact the DNA as the enzyme wraps around it. These residues could be varied to alter the binding of the protein. To obtain acceptable specificity, we may need to pick as initial DBP a mutant of EcoRI that folds and dimerizes, but that binds DNA weakly. The mutations in regions of the protein that contact DNA outside of the original GAATTC will confer the desired affinity and specificity on the novel protein.

[0160] One may wish to obtain a protein that binds to one target DNA sequence, but not to other sequences that contain a subsequence of the target. For example, we may seek a protein that recognizes TGAATTCA, but not any of the sequences vGAATTCb. To achieve this distinction, we place the target sequence in the promoter region of the selectable gene and one or more instances of the related sequences, to which we intend that the protein not bind, in the promoter region of an essential gene, such as an antibiotic-resistance gene.

[0161] Other stable proteins may also be used as initial DBPs, even if they show no DNA-binding properties. Par-  
raga *et al.* (Reference 8 in PARR88) report that Eisen *et al.* have fused 229 residues of yeast ADR1 to beta-galactosidase and that the fusion protein binds sequence-specifically to DNA *in vitro*.

[0162] Adenovirus E1A protein turns on early viral genes as well as the human heat shock protein hsp70 (SIMO88). Further, a normal inducible nuclear DNA-binding protein regulates the IL-2alpha interleukin-2 receptor-R(alpha) gene and also promotes activation of transcription from the HIV-1 virus LTR (BOHN88). These studies indicate one of the many difficulties of designing antiviral chemotherapy by using the transcriptional regulatory apparatus of the virus as a target. This invention uses unique target sequences, not represented elsewhere in the host genome, as targets for suppression of gene expression.

[0163] The DNA sequences of operators that interact with proteins that control mating-type and cell-type specific transcription in yeast (MILL85) reveal that the consensus site for action of the alpha2 protein dimer is symmetric, while a heterodimeric complex of alpha2 and a1 subunits acts on an asymmetric site. The alpha2a1-responsive site consists of a half-site that is identical to the alpha2 half-site, and another half-site that is a consensus for a1 protein binding. The spacings between the symmetric and asymmetric sites are not the same.

[0164] Antibodies that bind DNA and other nucleic acids have been obtained from human patients suffering from Systemic Lupus Erythematosus. Murine monoclonal antibodies have been obtained that specifically recognize Z-DNA, B-DNA, ssDNA, triplex DNA, and certain repeating sequences (ANDE88). Anderson *et al.* (ANDE88) report that: 1) the antibodies studied contact six base pairs and four phosphates, 2) antibodies are unlikely to provide some of the well known motifs for DNA-binding, e.g. helix-turn-helix, 3) study of DNA-antibody complexes may yield insights into mechanisms of recognition, and 4) a DNA-recognizing antibody might be converted into a sequence or structure specific nuclease. The shortness of the contact makes it unlikely that high specificity can be attained.

#### Properties of serially-linked globular domains:

[0165] A protein motif for DNA binding, present in some eukaryotic transcription factors, is the zinc finger in which zinc coordinately binds cysteine and histidine residues to form a conserved structure that is able to bind DNA (FRAN88). *Xenopus laevis* transcription factor TFIIIA is the first protein demonstrated to use this motif for DNA binding, but other proteins such as human transcription factor SP1, yeast transcription activation factor GAL4, and estrogen receptor protein have been shown to require zinc for DNA binding *in vitro* (EVAN88). Other mammalian and avian steroid hormone receptors and the adenovirus E1A protein, that bind DNA at specific sites, contain cysteine-rich regions which may form metal chelating loops.

[0166] Zinc-finger regions have been observed in the sequences of a number of eukaryotic DBPs, but no high-resolution 3D structure of a Zn-finger protein is yet available. A variety of models have been proposed for the binding of zinc-finger proteins to DNA (FAIR86, PARR88, BERG88, GIBS88). Model building suggests which residues in the Zn-fingers contact the DNA and these would provide the primary set of residues for variation. Berg (BERG88) and Gibson *et al.* (GIBS88) have presented models having many similarities but also some significant differences. Both models suggest that the motif comprises an antiparallel beta structure followed by an alpha helix and that the front side of the helix contacts the major groove of the DNA. By assuming that conserved basic residues of the Zn-finger make contact with phosphate groups in each copy of the motif, Gibson *et al.* deduce that the amino terminal part of the helix makes direct contact to the DNA. The Gibson model does not, however, account well for the number of bases contacted by Zn-finger proteins. The observations on H-T-H proteins suggest that a DNA-recognizing element can interact in a variety of ways with DNA and we assert that a similar situation is likely in Zn-finger proteins. Thus, until a 3D model of a Zn-finger protein bound to DNA is available, all of the residues modeled as occurring on the alpha helix away from the beta structure should be considered as primary candidates for variegation when one wishes to alter the DNA-binding properties of a Zn-finger protein. In addition, residues in the beta segment may control interactions with the sugar-phosphate backbone

which can effect both specific and non-specific binding.

[0167] Parraga *et al.* (PARR88) have reported a low-resolution structure of a single zinc-finger from NMR data. They confirm the alpha helix proposed by Berg and by Gibson *et al.*, but not the antiparallel beta sheet. The models proposed by Klug and colleagues (FAIR86) have a common feature that is at variance with the models of Berg and of Gibson *et al.*, viz. that the protein chain exits each finger domain at the same end that it entered. The structure published by Parraga *et al.* does not settle this point, but suggests that the exit strand tends toward the end opposite from the entrance strand, thereby supporting the overall models of Berg and of Gibson *et al.* Parraga *et al.* also report that a) a chimeric molecule consisting of zinc-finger domains linked to beta-galactosidase binds sequence-specifically to DNA and b) a protein comprising only two finger motifs can bind sequence-specifically to DNA. They do not suggest that the residues could be mutagenized to achieve novel recognition.

[0168] A protein composed of a series of zinc fingers offers the greatest potential of uniquely recognizing a single site in a large genome. A series of zinc fingers is not so well suited to development of a DBP that is sensitive to an effector molecule as is a more compact globular protein such as *E. coli* CRP. Positive control of genes adjacent to the target DNA subsequence can be achieved as in the case of TF-IIIa.

### Overview: Variegation Strategy

#### Choice of residues in parental potential-DBP to vary:

[0169] We choose residues in the initial potential-DBP to vary through consideration of several factors, including: a) the 3D structure of the initial DBP, b) sequences homologous to the initial DBP, c) modeling of the initial DBP and mutants of the initial DBP, d) models of the 3D structure of the target DNA, and e) models of the complex of the initial DBP with DNA. Residues may be varied for several reasons, including: a) to establish novel recognition by changing the residues involved directly in DNA contacts while keeping the protein structure approximately constant, b) to adjust the positions of the residues that contact DNA by altering the protein structure while keeping the DNA-contacting residues constant, c) to produce heterodimeric DBPs by altering residues in the dimerization interface while keeping DNA-contacting residues constant, and d) to produce pseudo-dimeric DBPs (see below) by varying the residues that join segments of dimeric DBPs while keeping the DNA-contacting residues and other residues fixed.

[0170] If a dimeric protein comprises two identical polypeptide chains related by a two-fold axis of rotation, we speak of a homodimer with two-fold dyad symmetry. When two very similar polypeptides fold into similar domains and associate, we may observe that there is an approximate two-fold rotational axis that relates homologous residues, such as the alpha1-beta1 dimer of haemoglobin. We refer to such a protein as a heterodimer and to the symmetry axis as a quasi-dyad. When we produce a single-chain DBP by fusing gene fragments that encode two DNA-binding domains joined by a linker amino acid subsequence, we call the molecule a pseudo-dimer and the axis that relates pairs of residues a pseudo-dyad.

#### Principles that guide choice of residues to vary:

[0171] A key concept is that only structured proteins exhibit specific binding, i.e. can bind to a particular chemical entity to the exclusion of most others. In the case of polypeptides, the structure may require stabilization in a complex with DNA. The residues to be varied are chosen to preserve the underlying initial DBP structure or to enhance the likelihood of favorable polypeptide-DNA interactions. The selection process eliminates cells carrying genes with mutations that prevent the DBP from folding. Genes that code for proteins or polypeptides that bind indiscriminately are eliminated since cells carrying such proteins are not viable. Although preservation of the basic underlying initial DBP structure is intended, small changes in the geometry of the structure can be tolerated. For example, the spatial relationship between the alpha 3 helix in one monomer of  $\lambda$  Cro and the alpha 3 helix in the dyad-related monomer (denoted alpha 3') is a candidate for variation. Small changes in the dimerization interface can lead to changes of up to several Å in the relative positions of residues in alpha 3 and alpha 3'.

[0172] Burial of hydrophobic surfaces so that bulk water is excluded is one of the strongest forces driving the folding of macromolecules and the binding of proteins to other molecules. Bulk water can be excluded from the region between two molecules or between two portions of a single molecule only if the surfaces are complementary. The double helix of B-DNA allows most of the hydrophobic surface nucleotides to be buried. The edges of the bases have several hydrogen-bonding groups; the methyl group of thymine is an important hydrophobic group in DNA (HARR88). To achieve tight binding, the shape of the protein must be highly complementary to the DNA, all or almost all hydrogen-bonding groups on both the DNA and the protein must make hydrogen bonds, and charged groups must contact either groups of opposite charge or groups of suitable polarity or polarizability.

[0173] There are two complementary interfaces of major interest: a) the DNA-protein interface and b) the interface between protein monomers of dimers or between domains of pseudo-dimers. The DNA-protein interface is more polar

than most protein-protein interfaces, but hydrophobic amino acids (e.g. F, L, M, V, I, W, Y) occur in sequence-specific DNA-protein interfaces. The protein-protein interfaces of natural DBPs are typical protein-protein interfaces.

[0174] Amino acids are classified as hydrophilic or hydrophobic (ROSE85, EISE86a,b), and although this classification is helpful in analyzing primary protein structures, it ignores that the side groups may contain both hydrophobic and hydrophilic portions, e.g., lysine. Hydrogen bonds and other ionic interactions have strong directional behavior, while hydrophobic interactions are not directional. Thus substitution of one hydrophobic side group for another hydrophobic side group of similar size in an interface is frequently tolerated and causes subtle changes in the interface. For the purposes of the present invention, such hydrophobic-interchange substitutions are made in the protein-protein interface of DBPs so that a) the geometry of the two monomers in the dimer will change, and b) compensating interactions produce exclusively heterodimers.

[0175] The process claimed here tests as many surfaces as possible to select one as efficiently as possible that binds to the target. The selection isolates cells producing those proteins that are more nearly complementary to the target DNA, or proteins in which intermolecular or intramolecular interfaces are more nearly complementary to each other so that the protein can fold into a structure that can bind DNA. The effective diversity of a variegated population is measured by the number of different surfaces, rather than the number of protein sequences. Thus we should maximize the number of surfaces generated in our population, rather than the number of protein sequences. Proteins do not have distinct, countable surfaces; therefore, we define an interaction set as a collection of residues of a protein that can simultaneously touch the target DNA.

[0176] If N spatially separated residues of a protein are varied,  $20 \times N$  surfaces are generated. Variation of N residues in the same interaction set yields  $20^N$  surfaces. For example, if  $N = 6$ , variation of spatially separated residues yields 120 surfaces while variation of interacting residues yields  $20^6 = 6.4 \times 10^7$  surfaces. The process of varying residues in an interaction set to maximize the number of surfaces obtained is referred to as Structure-directed Mutagenesis.

[0177] If the protein residues to be varied are close enough together in sequence that the variegated DNA (vgDNA) encoding all of them can be made in one piece, then cassette mutagenesis is picked. The present invention is not limited to a particular length of vgDNA that can be synthesized. With current technology, a stretch of 60 amino acids (180 DNA bases) can be spanned.

[0178] Mutation of residues further than sixty residues apart can be achieved using other methods, such as single-stranded-oligonucleotide-directed mutagenesis (BOTS85) and two or more mutating primers.

[0179] To vary residues separated by more than sixty residues, two cassettes may be mutated serially. From 2-fold to 1000-fold variegation is first introduced into a first cassette. We then introduce 1000-fold to  $10^6$ -fold variegation into a second cassette of the variegated vector population. The composite level of variation preferably does not exceed the prevailing capabilities to a) produce very large numbers of independently transformed cells or b) select small components in a highly varied population. The limits on the level of variegation are discussed below.

#### Assembly of Relevant Data:

[0180] Here we assemble the data about the initial DBP and the target that are useful in deciding which residues to vary in the variegation cycle:

- 1) 3D structure, or at least a list of residues that contact DNA and that are involved in the dimer contact of the initial DBP,
- 2) list of sequences homologous to the initial DBP, and
- 3) model of the target DNA sequence.

These data and an understanding of the function and structure of different amino acids in proteins will be used to answer three questions:

- 1) which residues of the initial DBP are on the outside and close enough together in space to touch the target DNA simultaneously?
- 2) which residues of the initial DBP can be varied with high probability of retaining the underlying initial DBP structure?
- 3) which residues of the initial DBP can affect the dimerization or folding of the initial DBP?



[0181] Although an atomic model of the target material is preferred in such examination, it is not necessary.

#### Graphical and computational tools:

5 [0182] The most appropriate method of picking the residues of the protein chain at which the amino acids should be varied is by viewing with interactive computer graphics a model of the initial DBP complexed with operator DNA. A model based on X-ray data from the DNA-protein complex is preferred, but other models may be used. A stick-figure representation of molecules is preferred. Suitable programs for viewing and manipulating protein and nucleic acid models include: a) PS-FRODO, written by T. A. Jones (JONE85) and distributed by the Biochemistry Department of Rice  
10 University, Houston, TX; and b) PROTEUS, developed by Dayringer, Tramantano, and Fletterick (DAYR86). Any hardware that supports either of these programs is appropriate.

#### Use of Knowledge of Mutations Affecting Protein Stability

15 [0183] In choosing the residues to vary and the substitutions to be made for such residues, one may make use not only of modelling as described above but also of experimental data concerning the effects of mutation in the initial DNA-binding protein. Mutations which will markedly reduce protein stability are to be avoided in most cases.

[0184] Missense mutations that decrease DNA-binding protein function non-specifically by affecting protein folding are distinguished from binding-specific mutations primarily on the basis of protein stability (NELS83, PAKU86, VERS86b, HECH84, HECH85a, and HECH85b).

20 [0185] Tables 1, 12, and 13 summarize the results of a number of studies on single missense mutations in the three bacteriophage repression proteins:  $\lambda$  repressor (Table 12) (NELS83, GUAR82, HECH85a, and NELS85),  $\lambda$  Cro (Table 1) (PAKU86, EISE85), and P22 Arc repressor (Table 13) (VERS86a, VERS86b). The majority of the mutant sequences shown in Tables 1, 12, and 13 were obtained in experiments designed to detect loss of function *in vivo*. The second-site  
25 pseudo-reversion mutations (HECH85a), and suppressed nonsense mutations (NELS83), restore function, and some of the site specific changes (EISE85) produce functional proteins.

[0186] Roughly 50-70% of the single missense mutations of the DNA-binding proteins selected for loss of function (Tables 1, 12, and 13) produce protein folding defects.

#### Use of Knowledge of Mutations Affecting the DNA-Protein Interface

[0187] Missense mutations in residues thought to be involved in specific interactions with DNA have been reported for several prokaryotic repressor proteins. Table 14 shows an alignment of the H-T-H DNA-binding domains of four  
35 prokaryotic repressor proteins (from top to bottom:  $\lambda$  repressor,  $\lambda$  Cro, 434 repressor and *trp* repressor) and indicates the positions of missense mutations in residues that are solvent-exposed in the free protein but become buried in the protein-DNA complex, and that affect DNA binding.

[0188] Randomly obtained missense mutations in solvent-exposed residues of  $\lambda$  repressor,  $\lambda$  Cro, and *trp* repressor, yield sets of mutants that reduce DNA binding (Table 14). These sets correlate well to the sets of residues that are proposed to interact directly with DNA. Some mutations in  $\lambda$  Cro (EISE85) and all those shown for 434 repressor  
40 (WHAR85a) were obtained through site-directed mutagenesis. Most of the mutations shown in the  $\lambda$  and *trp* repressor sequences are trans-dominant when the mutant gene is present on an overproducing plasmid (NELS83, KELL85). The exceptions to trans-dominance are the  $\lambda$  repressor SP35 and the *trp* repressor AT80 mutations. This latter change produces a repressor that has only slightly reduced binding (KELL85). The trans-dominance observed for these mutations is proposed by the authors to result from the wild-type repressor and the mutant repressor forming mixed oligomers  
45 which are inactive in binding to operator sites.

[0189] Wharton (WHAR85a) has reported that extensive site-directed mutagenesis of 434 repressor positions 28 and 29 produced no functional protein sequences other than the wild-type. Apparently, in the context of 434 repressor structure and operators, only proteins with the wild-type Q28-Q29 sequence bind to the wild-type operators.

[0190] Table 14 also shows missense mutations that result in near normal repressor activity. Substitution of 434  
50 repressor Q33 with H, L, V, T, or A produces repressors that function if expressed from overproducing plasmids (WHAR85a); repressor specificity is, however, reduced. Mutations in  $\lambda$  repressor, QY33 (NELS83, HECH83), and in  $\lambda$  Cro, YF26 (EISE85), produce altered proteins which make one less H-bond to the DNA and which bind to the operator DNA with reduced affinity. Thus, loss of a single H-bond is insufficient to completely abolish binding of DNA. Mutations YK26 and HR35 in  $\lambda$  Cro show nearly normal binding (EISE85).

55 [0191] Nelson and Sauer (NELS85) and Hecht *et al.* (HECH85a, b) have described four replacements in  $\lambda$  repressor (Table 12): EK34, GN48, GS48, and EK83. These derivatives have higher affinity for  $O_R1$  than w.t.  $\lambda$  repressor.

[0192] Extended amino acid arms at N- and C-terminal locations are important DNA-binding structures in at least four prokaryotic repressors:  $\lambda$  repressor and Cro, and P22 Arc and Mnt.

[0193] Sequence-specific and sequence-independent contacts are made by the first 6 amino acid residues (STKKKP) of the  $\lambda$  repressor N-terminal region which form an "arm" that can wrap around the DNA (ELIA85, PABO82a). Missense mutations KE4 and LP12 (Table 12) both greatly reduce repressor activity *in vivo* (NELS83). Deletion of the first six residues results in a protein which is non-functional *in vivo* (ELIA85). Deletion of the first three residues results in decrease of affinity for  $O_R1$ , loss of protection of back side guanines, altered specificity between  $O_R1$  and  $O_R3$ , and decreased binding sensitivity to changes in temperature or salt concentration (ELIA85, PABO82a).

[0194] Missense mutations of P22 Arc that produce non-functional proteins with high intracellular specific protein levels (Table 13) are found only in the N-terminal 10 residues of the protein (VERS86b). A single residue change at position 6 (HP6) in P22 Mnt changes operator recognition in the altered protein (YOUN83, VERS86a,b). Knight and Sauer (cited in VERS86a,b) replaced the first 6 residues of Mnt repressor with the first 9 residues of Arc repressor to produce a repressor that binds to the *arc* operator but not to the *mnt* operator. Thus P22 Mnt and Arc use a recognition region located in the first 6-10 amino-terminal residues for DNA recognition and binding. The N-terminal DNA-binding of these proteins can not be the recognition helix of a typical H-T-H motif.

[0195] In  $\lambda$  Cro, a C-terminal sequence (K62-K63-T64-T65-A66) has been suggested on the basis of model building (TAKE85) and NMR measurements (LEIG87) to form a flexible arm that interacts with minor groove phosphates. Eisenbeis and Caruthers (cited in KNIG88) have found that T64, T65, and A66 have minor effects on protein-operator affinity, while K63 is very important. The C-terminal sequence of P22 Mnt (K79-K80-T81-T82) is almost identical to that of  $\lambda$  Cro. It has been shown (KNIG88) that deletion of the three residues after K79 has little effect on protein structure or DNA binding. Deletion of K79 and the distal residues, however, reduces operator binding by three orders of magnitude with little apparent change in protein structure.

#### Use of Knowledge of Mutations Affecting the Protein-Protein Interface

[0196] It is also possible to modulate DNA-binding specificity by altering the protein-protein interface. Because the oligomerization equilibrium is coupled to DNA binding, mutations that alter oligomerization affect operator site affinity. Since oligomerization involves the matching of protein surfaces, many interactions are hydrophobic and mutations which specifically destabilize oligomerization are similar to mutations which destabilize global protein structure. Interactions at the site of oligomerization can influence the strength of interactions at the DNA-binding site by subtle alterations in protein structure.

#### Use of Mutations That Affect Activation

[0197] When  $\lambda$ , 434, and P22 repressors bind to their respective  $O_R2$  sites, they activate transcription (POTE80, POTE82, PTAS80). The site on  $\lambda$  repressor which activates RNA polymerase is located on the N-terminal domain of the molecule (BUSH88, HOCH83, SAUE79). Activation requires contact between the N-terminal domain of repressor at  $O_R2$  and RNA polymerase (HOCH83, SAUE79) and this contact stimulates isomerization of the polymerase complex to the open form (McClure and Hawley, cited in GUAR82).

[0198] Missense mutations in  $\lambda$ , P22, or 434 repressors that specifically reduce  $P_{RM}$  activation while leaving operator binding intact are in the solvent-exposed protein surface closest to RNA polymerase bound at  $P_{RM}$  (GUAR82, PABO79, BUSH88, WHAR85a). For  $\lambda$  and 434 repressor this surface includes residues in alpha helix 2 and in the turn between alpha helices 2 and 3. In P22 repressor, the surface is formed at the carboxyl terminus of alpha helix 3 (PABO79, TAKE83). In each repressor, the changes that reduce transcriptional activation at  $P_{RM}$  involve the substitution of a basic residue for a neutral or acid residue. Further, missense mutations in  $\lambda$  and 434 repressors which increase transcription at  $P_{RM}$  involve the substitution of an acidic residue for a neutral or basic residue (GUAR82, BUSH88).

[0199] Transcriptional activation at  $P_{RM}$  involves the apposition of a negatively charged surface on the N-terminal domain of  $\lambda$ , 434, or P22 repressor to a site on RNA polymerase (BUSH88). Mutations that a) alter the negatively-charged surface of repressor by removing acidic residues or by replacing them with basic residues, or b) that position the negative surface incorrectly with respect to RNA polymerase, decrease transcriptional activation at  $P_{RM}$ . Alterations that produce a more negatively charged surface act to increase transcription at  $P_{RM}$ .

#### Pick principal set of residues to vary:

[0200] A huge number of variant DNA sequences can be generated by synthesis with mixed reagents at chosen bases. Usually, it is necessary that the number of variants not exceed the number of independently transformed cells generated from the synthetic DNA. It is efficient, however, to make the number of variants as close as practical to this limit. The total number of variants is the product of the number of variants at each varied codon over all the variable codons. Thus, we first consider which residues could be varied with an expectation that alteration could affect DNA binding. We then pick a range of amino acids at each variable residue. The total number of variants is the product of

these numbers. If the product is too large or too small, we alter the list of residues and range of variation at each variable residue until an acceptable number is found.

[0201] Considering which residues are on the surface of the initial DBP, we pick residues that are close enough together on the surface of the initial DBP to touch a molecule of the target simultaneously without having any initial DBP main-chain atom come closer than van der Waals distance (viz. 4.0 to 5.0 Å center to center) to any target atom. For the purposes of the present invention, a residue of the initial DBP "touches" the target if:

- a) a main-chain atom is within van der Waals distance, viz. 4.0 to 5.0 Å, of any atom of the target molecule,
- b) the C<sub>beta</sub> is within a specific distance of any atom of the target molecule so that a side-group atom could make contact with that atom, or
- c) there is evidence that altering the residue alters the DNA-binding of the initial DBP.

[0202] The residues in the principal set need not be contiguous in the protein sequence. The exposed surfaces of the residues to be varied need not be connected. We prefer only that the amino acids in the residues to be varied all be capable of touching a single copy of the target DNA sequence simultaneously without atoms overlapping.

[0203] In addition to the geometrical criteria, we prefer that there be indications that the initial DBP structure will tolerate substitutions at each residue in the principal set of residues. Indications could come from various sources, including homologous sequences and modeling.

#### Pick a secondary set of residues to vary:

[0204] The secondary set comprises those residues not in the primary set that touch residues in the primary set. These residues might be excluded from the primary set because the residue is : a) internal, b) highly conserved, or c) on the surface, but the curvature of the initial DBP surface prevents the residue from being in contact with the target at the same time as one or more residues in the primary set.

[0205] Internal residues are frequently conserved and the amino acid type can not be changed to a significantly different type without risk that the protein structure will be disrupted. Nevertheless, some conservative changes of internal residues, such as I to L or F to Y, are tolerated. Such conservative changes affect the detailed placement and dynamics of adjacent protein residues and such variation may be useful to improve the characteristics of DBP binding.

[0206] Surface residues in the secondary set are most often located on the periphery of the principal set. Such peripheral residues can not make direct contact with the target simultaneously with all the other residues of the principal set. It is appropriate to vary the charge of some or all of these residues. For example, the variegated codon containing equimolar A and G at base 1, equimolar C and A at base 2, and A at base 3 yields amino acids T, A, K, and E with equal probability.

#### Choice of residues to vary simultaneously:

[0207] The allowed level of variegation determines how many residues can be varied at once; geometry determines which ones. The user may pick residues to vary in many ways; the following is a preferred manner. The user picks the objective of the variegation, vide supra.

[0208] The number of residues picked is coupled to the range through which each can be varied. In the first round progressivity is not an issue; the user may elect to produce a level of variegation such that each molecule of vgDNA is potentially different through, for example, unlimited variegation of 10 codons ( $20^{10}$  approx. =  $10^{13}$  different protein sequences). The levels of efficiency of ligation and transformation reduce the number of DNA sequences actually tested to between  $10^7$  and  $10^9$ . Multiple performances of the process with very high levels of variegation will not yield repeatable results; the user decides whether this is important.

#### Pick range of variation:

[0209] Each varied residue can have a different scheme of variegation, producing 2 to 20 different possibilities. We require that the process be progressive, i.e. each variegation cycle produces a better starting point for the next variegation cycle than the previous cycle produced.

N.B.: Setting the level of variegation such that the parental pdbp and many sequences related to the parental pdbp sequence are present in detectable amounts insures that the process is progressive. If the level of variegation is so high that the frequency of the parental pdbp sequence can not be detected as a transformant, then each round of mutagenesis is independent of previous rounds and there is no assurance of progressivity. This approach can lead to valuable DNA-binding proteins, but multiple repetitions of the process at this level of variegation will not yield pro-

gressive results. Excessive variegation is not preferred in subsequent iterations of this process.

[0210] Progressivity is not an all-or-nothing property. So long as most of the information obtained from previous variegation cycles is retained and many different surfaces that are related to the parental DBP surface are produced, the process is progressive. If the level of variegation is so high that the parental *dbp* gene may not be detected, the assurance of progressivity diminishes. If the probability of recovering the parental DBP is negligible, then the probability of progressive results is also negligible.

[0211] An opposing force in our design considerations is that DBPs are useful in the population only up to the amount that can be detected; any excess above the detectable amount is wasted. Thus we produce as many surfaces related to the parental DBP as possible within the constraint that the parental DBP be present as a marker for the detection level.

#### Mutagenesis of DNA:

[0212] We now decide how to distribute the variegation within the codons for the residues to be varied. These decisions are influenced by the nature of the genetic code. When vgDNA is synthesized, variation at the first base of a codon creates a population coding for amino acids from the same column of the genetic code table (Table 16); variation at the second base of the codon creates a population coding for amino acids from the same row of the genetic code table; variation at the third base of the codon creates a population coding for amino acids from the same box. Work with 3D protein structural models may suggest definite sets of amino acids to substitute at a given residue, but the method of variation may require either more or fewer kinds of amino acids be included. For example, substitution of N or Q at a given residue may be wanted. Combinatorial variation of codons requires that mixing N and Q at one location also include K and H as possibilities at the same residue. The present invention does not rely on accurate predictions of the amino acids to be placed at each residue, rather attention is focused on which residues should be varied.

[0213] There are many ways to generate diversity in a protein (RICH86, CARU85, OLIP86). An extreme case is that one or a few residues of the protein are varied as much as possible (*inter alia* see CARU85, CARU87, RICH86, WHAR85a). We will call this limit "Focused Mutagenesis". When there is no binding between the parental DBP and the target, we preferably pick a set of five to seven residues on the surface and vary each through all 20 possibilities.

[0214] An alternative plan of mutagenesis ("Diffuse Mutagenesis") that may be useful is to vary many more residues through a more limited set of choices (VERS86a,b, INOU86 (Ch.15), PAKU86). This can be accomplished by spiking each of the pure nucleotides activated for DNA synthesis (e.g. nucleotide-phosphoramidites) with one or more of the other activated nucleotides. Contrary to general practice, the present invention sets the level of spiking so that only a small percentage (1% to .00001%, for example) of the final product will contain the parental DNA sequence. This will insure that the majority of molecules carry single, double, triple, and higher mutations and, as required for progressivity, that recovery of the parental sequence will be a possible outcome.

[0215] Let  $N_b$  be the number of bases to be varied, and let  $Q$  be the fraction of all DNA sequences that should have the parental sequence, then  $M$ , the fraction of the nucleotide mixture that is the majority component, is

$$M = \exp\{\log_e(Q)/N_b\} = 10(\log_{10}(Q)/N_b).$$

If, for example, thirty base pairs on the DNA chain were to be varied and 1% of the product is to have the parental sequence, then each mixed nucleotide substrate should contain 86% of the parental nucleotide and 14% of other nucleotides. Table 17 shows the fraction ( $f_n$ ) of DNA molecules having  $n$  non-parental bases when 30 bases are synthesized with reagents that contain fraction  $M$  of the majority component. When  $M=.63096$ ,  $f_{24}$  and higher are less than  $10^{-8}$ . Note that substantial probability for 8 or more substitutions occurs only if the fraction of parental sequence ( $f_0$ ) drops to around  $10^{-3}$ .

[0216] The  $N_b$  base pairs of the DNA chain that are synthesized with mixed reagents need not be contiguous. They are picked so that between  $N_b/3$  and  $N_b$  codons are affected to various degrees. The residues picked for mutation are picked with reference to the 3D structure of the initial DBP, if known. For example, one might pick all or most of the residues in the principal and secondary set. We may impose restrictions on the extent of variation at each of these residues based on homologous sequences or other data. The mixture of non-parental nucleotides need not be random, rather mixtures can be biased to give particular amino acid types specific probabilities of appearance at each codon. For example, one residue may contain a hydrophobic amino acid in all known homologous sequences; in such a case, the first and third base of that codon would be varied, but the second would be set to T. This Diffuse Mutagenesis will reveal the subtle changes possible in the protein backbone associated with conservative interior changes, such as V to I, as well as some not so subtle changes that require concomitant changes at two or more residues of the protein.

Focused Mutagenesis:

[0217] If we have no information indicating that a particular amino acid or class of amino acid is appropriate, we approximate substitution of all amino acids with equal probability because representation of one or a few pdbp genes above the detectable level is unproductive. Equal amounts of all four nucleotides at each position in a codon yields the amino acid distribution:

|        |        |           |        |        |        |
|--------|--------|-----------|--------|--------|--------|
| 4/64 A | 2/64 C | 2/64 D    | 2/64 E | 2/64 F | 4/64 G |
| 2/64 H | 3/64 I | 2/64 K    | 6/64 L | 1/64 M | 2/64 N |
| 4/64 P | 2/64 Q | 6/64 R    | 6/64 S | 4/64 T | 4/64 V |
| 1/64 W | 2/64 Y | 3/64 stop |        |        |        |

[0218] This distribution has the disadvantage of giving two basic residues for every acidic residue. Such predominance of basic residues is likely to promote sequence-independent DNA binding. In addition, six times as much R, S, and L as W or M occur for the random distribution. Use of equimolar C and G at the third base reduces the over-representation of S, R, and L, but does not cure the maldistribution of acidics and basics.

[0219] Consider the distribution of amino acids encoded by one codon in a population of vgDNA. Let  $Abun(x)$  be the abundance of DNA sequences coding for amino acid  $x$ . For any distribution, there will be a most-favored amino acid (mfaa) with abundance  $Abun(mfaa)$  and a least-favored amino acid (lfaa) with abundance  $Abun(lfaa)$ . We seek the nucleotide distribution that allows all twenty amino acids and that yields the largest ratio  $Abun(lfaa)/Abun(mfaa)$  subject to two constraints. First, the abundances of acidic and basic amino acids should be equal. Second, the number of stop codons should be kept as low as possible. Thus only nucleotide distributions that yield

$$Abun(E)+Abun(D) = Abun(R)+Abun(K)$$

are considered, and the function maximized is:

$$f(\text{distribution}) = \{(1 - Abun(\text{stop})) (Abun(lfaa)/Abun(mfaa))\}.$$

[0220] We limit the third base to equimolar T and G (C and G would be equivalent). All amino acids are possible and the number of accessible stop codons is reduced.

[0221] A computer program, "Find Optimum vgCodon." (Table 18), varies the composition at bases 1 and 2, in steps of 0.05, and reports the composition that gives the largest value of  $f(\text{distribution})$  subject to the constraints:

$$g2 = (g1 \cdot a2 - 0.5 \cdot a1 \cdot a2) / (c1 + 0.5 \cdot a1),$$

$$t1 = 1 - a1 - c1 - g1, \text{ and}$$

$$t2 = 1 - a2 - c2 - g2.$$

The first constraint requires equal amount of acidic and basic amino acids and the second and third conserve matter. We vary  $a1$ ,  $c1$ ,  $g1$ ,  $a2$ , and  $c2$  and then calculate  $t1$ ,  $g2$ , and  $t2$ . Initially, variation is in steps of 5%. Once an approximately optimum distribution of nucleotides is determined, the region is further explored with steps of 1%. The optimum distribution is:

Optimum vgCodon

[0222]

|           | T    | C    | A    | G    |
|-----------|------|------|------|------|
| base #1 = | 0.26 | 0.18 | 0.26 | 0.30 |

(continued)

|           | T    | C    | A    | G    |
|-----------|------|------|------|------|
| base #2 = | 0.22 | 0.16 | 0.40 | 0.22 |
| base #3 = | 0.5  | 0.0  | 0.0  | 0.5  |

and yields DNA molecules encoding each type of amino acid with the abundances shown in Table 19.

[0223] The actual nucleotide distribution obtained in synthetic DNA will differ from the specified nucleotide distribution due to several causes, including: a) differential inherent reactivity of nucleotide substrates, and b) differential deterioration of reagents. It is possible to compensate partially for these effects, but some residual error will occur. We denote the average discrepancy between specified and observed nucleotide fraction as  $S_{err}$ .

$$S_{err} = \text{square root ( average[ } (f_{obs} - f_{spec})/f_{spec} \text{ ] )}$$

where  $f_{obs}$  is the amount of one type of nucleotide found at a base and  $f_{spec}$  is the amount of that type of nucleotide that was specified at the same base. The average is over all specified types of nucleotides and over a number (e.g. 10 to 50) of different variegated bases. By hypothesis, the actual nucleotide distribution at a variegated base will be within 5% of the specified distribution. Actual DNA synthesizers and DNA synthetic chemistry may have different error levels. It is the user's responsibility to determine  $S_{err}$  for the DNA synthesizer and chemistry employed by the user.

[0224] To determine the possible effects of errors in nucleotide composition on the amino acid distribution, we modified the program "Find Optimum vgCodon" in four ways:

1) the fraction of each nucleotide in the first two bases is allowed to vary from its optimum value times  $(1 - S_{err})$  to the optimum value times  $(1 + S_{err})$  in seven equal steps ( $S_{err}$  is the hypothetical fractional error level), maintaining the sum of nucleotide fractions for one codon position at 1.0.

2) g2 is varied in the same manner as a2, i.e. we dropped the restriction that  $Abun(D) + Abun(E) = Abun(K) + Abun(R)$ .

3) t3 and g3 are varied from 0.5 times  $(1 - S_{err})$  to 0.5 times  $(1 + S_{err})$  in three equal steps.

4) the smallest ratio  $Abun(lfaa)/Abun(mfaa)$  is sought.

In actual experiments, we direct the synthesizer to produce the optimum DNA distribution "Optimum vgCodon" given above. Incomplete control over DNA chemistry may, however, cause us to actually obtain the following distribution that is the worst that can be obtained if all nucleotide fractions are within 5% of the amounts specified in "Optimum vgCodon". A corresponding table can be calculated for any given  $S_{err}$  using the program "Find worst vgCodon within  $S_{err}$  of given distribution." given in Table 20.

#### Optimum vgCodon, worst 5% errors

[0225]

|           | T     | C     | A     | G     |
|-----------|-------|-------|-------|-------|
| base #1 = | 0.251 | 0.189 | 0.273 | 0.287 |
| base #2 = | 0.209 | 0.160 | 0.400 | 0.231 |
| base #3 = | 0.475 | 0.0   | 0.0   | 0.525 |

[0226] This distribution yields DNA encoding each of the twenty amino acids at the abundances shown in Table 21.

[0227] Each codon synthesized with the distribution of bases shown above displays  $4 \times 4 \times 2 = 2^5 = 32$  possible DNA sequences, though not in equal abundances. An oligonucleotide containing N such codons would display  $2^{5N}$  possible DNA sequences and would encode  $20^N$  protein sequences. Other variegation schemes produce different numbers of DNA and protein sequences. For example, if two bases in one codon are varied through two possibilities each, then

there are  $2 \times 2 = 4$  DNA sequences and  $2 \times 2 = 4$  protein sequences.

[0228] If five codons are synthesized with reagents mixed so as to produce the nucleotide distribution "Optimum vgCodon", and if we actually obtained the nucleotide distribution "Optimum vgCodon, worst 5% errors", then DNA sequences encoding the mfaa at all of the five codons are about 277 times as likely as DNA sequences encoding the lfaa at all of the five codons. Further, about 24% of the DNA sequences will have a stop codon in one or more of the five codons.

[0229] Consider variegation of a hypothetical sequence, F24-G25-D26-E27-T28, in which each variegated codon is synthesized as an "Optimal vgCodon". The actual abundance of the DNA encoding each type of amino acid is, however, taken from the case of  $S_{err} = 5\%$  given in Table 21. The abundance of DNA encoding the parental amino acid sequence is:

**Amount (parental seq.)**

$$\begin{aligned} & \text{F24} \quad \text{G25} \quad \text{D26} \quad \text{E27} \quad \text{T28} \\ = & \text{Abun(F)} * \text{Abun(G)} * \text{Abun(D)} * \text{Abun(E)} * \text{Abun(T)} \\ = & .0249 \times .0663 \times .0545 \times .0602 \times .0437 \\ = & 2.4 \times 10^{-7} \end{aligned}$$

Therefore, if the efficiency of the entire process allows us to examine  $10^7$  different DNA sequences, DNA encoding the parental DBP sequence as well as very many related sequences will be present in sufficient quantity to be detected and we are assured that the process will be progressive.

#### Setting level of variegation:

[0230] We use the following procedure to determine whether a given level of variegation is practical:

1) from: a) the intended nucleotide distribution at each base of a variegated codon, and b)  $S_{err}$  (the error level in mixed DNA synthesis), calculate the abundances of DNA sequences coding for each amino acid and stop,

2) calculate the abundance of DNA encoding the parental DBP sequence by multiplying the abundances of the parental amino acid at each variegated residue,

The abundances used in the procedure above are calculated from the worst distribution that is within  $S_{err}$  of the specified distribution. A variegation that insures that the parental DBP sequence can be recovered is practical. Such a level of variegation produces an enormous number of multiple changes related to the parental DBP available for selection of improved successful DBPs. We adjust the subset of residues to be varied and levels of variegation at each residue until the calculated variegation is within bounds.

#### Reduction of gratuitous restriction sites:

[0231] If the method of mutagenesis to be used is replacement of a cassette, we consider whether the variegation generates gratuitous restriction sites. We reduce or eliminate gratuitous restriction sites by appropriate choice of variegation pattern and silent alteration of codons neighboring the sites of variegation.

#### Focused mutagenesis:

[0232] In the preferred embodiment of this process, the number of residues and the range of variation at each residue are chosen to maximize the number of DNA binding surfaces, to minimize gratuitous restriction sites, and to assure the recovery of the initial DBP sequence. For example, in Detailed Example 1, the initial DBP is  $\lambda$  Cro. One primary set of residues includes G15, Q16, K21, Y26, Q27, S28, N31, K32, H35, A36, and R38 of the H-T-H region (Table 14b) and C-terminal residues K56, N61, K62, K63, T64, T65, and A66. A secondary set of residues includes L23, G24, and V25 from the turn portion of the H-T-H region, buried residues T20, A21, A30, I31, A34, and I35 from alpha helices 2 and 3, and dimerization region residues E54, V55, F58, P59, and S60.

[0233] The initial set of 5 residues for Focused Mutagenesis contains residues in or near the N-terminal half of alpha helix 3: Y26, Q27, S28, N31, and K32. Varying these 5 residues through all 20 amino acids produces  $3.2 \times 10^6$  different protein sequences encoded by  $32^5 (=3.3 \times 10^7)$  different DNA sequences. Since all 5 residues are in the same

interaction set, this variegation scheme produces the maximum number of different surfaces. Assuming optimized nucleotide distribution described above and  $S_{\text{err}} = 5\%$ , the probability of obtaining the parental sequence is  $3.2 \times 10^{-7}$ . This level is within bounds for synthesis, ligation, transformation, and selection capable of examining  $10^8$  sequences of vgDNA. Codons for the 5 residues picked for Focused Mutagenesis are contained in the 51 bp PpuMII to BglII fragment of the ray<sup>+</sup> gen constructed in Detailed Example 1.

#### Repetition to obtain desired degree of DNA-binding:

[0234] The first variegation step can produce one or more DBPs having DNA-binding properties that are satisfactory to the user. If the best selected DBP is not fully satisfactory, parental DBPs for a second variegation step are picked from DBPs isolated in the first variegation step. The second and subsequent variegation steps may employ either Focused or Diffuse Mutagenesis procedures on residues of the primary or secondary sets. In the preferred embodiment of this process, the user chooses residues and mutagenesis procedures based on the structure of the parental DBP and specific goals. For example, consider three hypothetical cases.

[0235] In a first case, a variegation step produces a DBP with greater non-specific DNA binding than is desired. Information from sequence analysis and modeling is used to identify residues involved in sequence independent interactions of the DBP with DNA in the non-specific complex. In the next variegation step, some or all of these residues, together with one or more additional residues from the primary set, are chosen for Focused Mutagenesis and additional residues from the primary or secondary sets are chosen for Diffuse Mutagenesis.

[0236] In a second hypothetical case, a variegation step produces a DBP with strong sequence specific binding to the target and the goal is to optimize binding. In this case, the next variegation step employs Diffuse Mutagenesis of a large number of residues chosen mostly from the secondary set.

[0237] In the third hypothetical case, a DBP has been isolated that has insufficient binding properties. A set of residues is chosen to include some primary residues that have not been subjected to variation, one or more primary residues that have been varied previously, and one or more secondary residues. Focused Mutagenesis is performed on this set in the next variegation step.

#### Overview: DNA Synthesis, Purification, and Cloning

##### DNA sequence design:

[0238] The present invention is not limited to a single method of gene design. The idbp gene need not be synthesized in toto; parts of the gene may be obtained from nature. One may use any genetic engineering method to produce the correct gene fusion, so long as one can easily and accurately direct mutations to specific sites. In all of the methods of mutagenesis considered in the present invention, however, it is necessary that the DNA sequence for the idbp gene be unique compared to other DNA in the operative cloning vector. If the method of mutagenesis is to be replacement of subsequences coding for the potential-DBP with vgDNA, then the subsequences to be mutagenized must be bounded by restriction sites that are unique with respect to the rest of the vector. If single-stranded oligonucleotide-directed mutagenesis is to be used, then the DNA sequence of the subsequence coding for the initial DBP must be unique with respect to the rest of the vector.

[0239] The coding portions of genes to be synthesized are designed at the protein level and then encoded in DNA. The amino acid sequences are chosen to achieve various goals, including: a) expression of initial DBP intracellularly, and b) generation of a population of potential-DBPs from which to select a successful DBP. The ambiguity in the genetic code is exploited to allow optimal placement of restriction sites and to create various distributions of amino acids at variegated codons.

##### Organization of gene synthesis:

[0240] The present invention is not limited as to how a designed DNA sequence is divided for easy synthesis. An established method is to synthesize both strands of the entire gene in overlapping segments of 20 to 50 nucleotides (THER88). An alternative method that is more suitable for synthesis of vgDNA is similar to methods published by others (OLIP86, OLIP87, AUSU87, KARN84). Contrary to most previous workers, we: a) use two synthetic strands, and b) do not cut the extended DNA in the middle. Our goals are: a) to produce longer pieces of dsDNA than can be synthesized as ssDNA on commercial DNA synthesizers, and b) to produce strands complementary to single-stranded vgDNA. By using two synthetic strands, we remove the requirement for a palindromic sequence at the 3' end. Moreover, the overlap should not be palindromic lest single DNA molecules prime themselves.

[0241] The present invention is not limited to any particular method of DNA synthesis or construction. Preferably, DNA is synthesized on a Milligen 7500 DNA synthesizer (Milligen, Bedford, MA) by standard procedures. Synthetic



DNA is purified by polyacrylamide gel electrophoresis (PAGE) or high-pressure liquid chromatography (HPLC). The present invention is not limited to any particular method of purifying DNA for genetic engineering.

#### IDBP Gene cloning:

[0242] We clone the *idbp* gene using plasmids that are transformed into competent bacterial cells by standard methods (MANI82) or slightly modified standard methods. DNA fragments derived from nature are operably linked to other fragments of DNA.

[0243] Cells transformed with the plasmid bearing the complete *idbp* gene are tested to verify expression of the initial DBP. Selection for plasmid presence is maintained on all media, while selections for DBP<sup>+</sup> phenotypes are applied only after growth in the presence of inducer appropriate to the promoter. Colonies that display the DBP<sup>+</sup> phenotypes in the presence of inducer and DBP<sup>-</sup> phenotypes in the absence of inducer are retained for further genetic and biochemical characterization. The presence of the *idbp* gene is initially detected by restriction enzyme digestion patterns characteristic of that gene and is confirmed by sequencing.

[0244] The dependence of the IDBP<sup>+</sup> and IDBP<sup>-</sup> phenotypes on the presence of this gene is demonstrated by additional genetic constructions. These are a) excision of the *idbp* gene by restriction digestion and closure by ligation, and b) ligation of the excised *idbp* gene into a plasmid recipient carrying different markers and no *idbp* gene. Plasmids obtained by excising the gene confer the DBP<sup>-</sup> phenotypes (e.g. Tc<sup>R</sup>, Fus<sup>S</sup>, and Gal<sup>S</sup> in Detailed Example 1). Plasmids obtained from ligation of *idbp* to a recipient plasmid confer the DBP<sup>+</sup> phenotypes in the presence of an inducer appropriate to the regulatable promoter (e.g. Tc<sup>S</sup>, Fus<sup>R</sup>, and Gal<sup>R</sup> in Detailed Example 1). Finally, a most important demonstration of the successful construction involves determination of the quantitative dependence of the selected phenotypes on the exogenous inducer concentration.

#### Overview: DNA-binding Protein Purification and Characterization

##### Isolation of IDBP:

[0245] We purify IDBP and its derivatives by standard methods, such as those described in JOHN80, TAKE86, LEIG87, VERS85b, KADO86.

##### Quantitation and characterization of protein-DNA binding:

[0246] Methods that can be used to quantitate and characterize sequence-specific and sequence-independent binding of a DBP to DNA include: a) filter-binding assays, b) electrophoretic mobility shift analysis, and c) DNase protection experiments. Ionic strength, pH, and temperature are important factors influencing DBP binding to DNA. Standard conditions should correspond closely to the anticipated conditions of use. Thus, if a binding protein is intended for use in bacterial cells in standard culture, a reasonable range of values from which to choose standard conditions would be: pH=7.5 to 8.0, 0.1 to 0.2 M KCl, and 32° to 37°C. Assay buffers preferably include cofactors, stabilizing agents, and counter ions for proper DBP function.

[0247] We prepare DNA fragments for analysis of protein-DNA binding by methods that are very similar to those described in MAXA77, KLEN70, RIGB77, and KIMJ87. Filter-binding assays can yield thermodynamic ( $K_D$ ) and kinetic ( $k_a$  and  $k_d$ ) constants and are performed by methods similar to those described by RIGG70, and KIMJ87. Electrophoretic mobility shift measurements can also yield values of  $K_D$ ,  $k_a$ , and  $k_d$  and are performed by methods similar to those of FRIE81. DNase protection assays use the methods of JOHN79, MAXA77, FOXK88. We use chemical methods to characterize binding of proteins to DNA similar to the methods described in BRUN87, BUSH85, and JENJ86.

Table of Examples

|             |  |
|-------------|--|
| <u>Ex.1</u> | Protocol for developing a new DNA-binding protein with affinity for a DNA-sequence found in HIV-1, by variegation of $\lambda$ Cro.  |
| <u>Ex.2</u> | Protocol for developing a new DNA-binding polypeptide with affinity for a DNA-sequence found in HIV-1, by variegation of a polypeptide having a segment homologous with Phage P22 Arc. |
| <u>Ex.3</u> | Use of a custodial domain (residues 20-83 of barley chymotrypsin inhibitor) to protect a DNA-binding polypeptide from degradation.   |

(continued)

| Table of Examples |   |
|-------------------|---|
| Ex. 4             | Use of a custodial domain containing a DNA-recognizing element (alpha-3 helix of Cro) to protect a DNA-binding polypeptide from degradation.  |
| Ex. 5             | Protocol for addition of arm to Phage P22 ARc to alter its DNA-binding characteristics.   |
| Ex. 6             | Protocol for preparation of novel DNA-binding protein that recognizes an asymmetric DNA sequence and corresponds to a fusion of third zinc-finger domain of the <i>Drosophila</i> <i>kr</i> gene product and the DNA-binding domain of Phage P22 Arc. |

**DETAILED EXAMPLE 1**

[0248] Below is a hypothetical example of a protocol for developing a new DNA-binding protein derived from  $\lambda$  Cro with affinity for a DNA sequence found in human immunodeficiency virus type 1 (HIV-1) using *E. coli* K-12 as the cell line or strain. Further optimization, in accordance with the teachings herein, may be necessary to obtain the desired results. Possible modifications in the preferred method are discussed following various steps of the example.

[0249] By hypothesis, we set the following technical capabilities:

## Yield from DNA synthesis

---

500 ng/synthesis of ssDNA 100 bases long,

10 ug/synthesis of ssDNA 60 bases long,

1 mg/synthesis of ssDNA 20 bases long.

---

## Maximum oligonucleotide

---

100 bases

---

## Yield of plasmid DNA

---

1 mg/l of culture medium

---

## Efficiency of DNA Ligation

---

0.1 % for blunt-blunt,

4 % for sticky-blunt,

11 % for sticky-sticky.

---

## Yield of transformants

---

$5 \times 10^8$  / ug DNA

---

Error in mixed DNA synthesis ( $S_{err}$ )

---

5%

---

Choice of cell line or strain:

[0250] In this example, the following *E. coli* K-12 *recA* strains are used: ATCC #35,882 *delta4* (Genotype: W3110 *trpC*, *recA*, *rpsL*, *sup<sup>o</sup>*, *delta4* (*gal-chlD-pql-att<sub>lambda</sub>*)) and ATCC #33,694 HB101 (Genotype: F', *leuB*, *proA*, *recA*, *thi*, *ara*, *lacY*, *galK*, *xyl*, *mtl*, *rpsL*, *supE*, *hsdS*, (*rB*, *mB*)). *E. coli* K-12 strains are grown at 37°C in LB broth (MANI82, p440) and on LB agar (addition of 15 g Bacto-agar) for routine purposes. Selections for plasmid uptake and maintenance are performed with addition of ampicillin (Ap) (200 ug/ml), tetracycline (Tc) (12.5 ug/ml) and kanamycin (Km) (50 ug/ml).

Choice of initial DBP:

[0251] The initial DBP is  $\lambda$  Cro. Helix-turn-helix proteins are preferred over other known DBPs because more detail is known about the interactions of these proteins with DNA than is known for other classes of natural DBP.  $\lambda$  Cro is pre-

ferred over  $\lambda$  repressor because it has lower molecular weight. Cro from 434 is smaller than  $\lambda$  Cro, but more is known about the genetics and 3D structure of  $\lambda$  Cro. An X-ray structure of the  $\lambda$  Cro protein has been published, but no X-ray structure of a DNA-Cro complex has appeared. A mutant of  $\lambda$  Cro, Cro67, confers the positive control phenotype *in vitro* but not *in vivo*. The contacts that stabilize the Cro dimer are known, and several mutations in the dimerization function have been identified (PAKU86).

[0252] By the methods disclosed herein, DBPs may be developed from Cro which recognize DNA binding sites different from the  $\lambda$  O<sub>P3</sub> or  $\lambda$  operator consensus binding sites, including heterodimeric DBPs which recognize non-symmetric DNA binding sites.

#### Selections for phenotypes conferred by DBP<sup>+</sup> function:

[0253] Media generally are supplemented with IPTG and antibiotic for selection of plasmid maintenance. Cell background is generally strain delta4 (galK.T.E deletion).

a. Galactose resistance (Gal<sup>R</sup>). Galactose epimerase deficient (galE<sup>-</sup>) strains of *E. coli* (BUTT63) lyse when treated with galactose. Selective medium is supplemented with 2% galactose, added after autoclaving. Additional galactose, up to 8%, somewhat reduces the background of artifactual galactose-sensitive colonies.

b. Galactose resistance selected immediately after transformation. Inducer IPTG is added to transformed cells, to  $5 \times 10^{-4}$  M at the start of the growth period, that allows expression of plasmid antibiotic-resistance. At 60 min after heat shock, cells are further diluted 10-fold into fresh LB broth containing IPTG, antibiotic to select for plasmid uptake (e.g. Ap or Km), and 2% galactose. Cells are grown until lysis is complete or for 3 h, whichever occurs first, then centrifuged at 6,000 rpm for 10 min, resuspended in the initial volume of the post-transformation growth culture, and applied to medium for further selection.

c. Fusaric acid resistance (Tc<sup>S</sup>, Fus<sup>R</sup>). Successful repression of tet yields resistance to lipophilic chelating agents such as fusaric acid (Fus<sup>R</sup> phenotype). Medium described by MALO81 is used for selection of fusaric acid resistance in *E. coli*; the amount of fusaric acid may be varied. Total cell inoculum is not greater than  $5 \times 10^6$  per plate.

d. Fusaric acid resistance and galactose resistance. Galactose at a final concentration of 2% is added to the medium described by MALO81 after autoclaving. Cells selected directly for galactose resistance in liquid following transformation are applied to this medium.

#### Selections for phenotypes conferred by DBP<sup>-</sup> function:

[0254] Cell background is generally strain HB101 (galK<sup>-</sup>). Media are generally supplemented with IPTG and antibiotic for plasmid maintenance.

a. Tc resistance. Medium, usually LB agar, is supplemented with Tc after autoclaving. Tc stock solution is 12.5 mg/ml in ethanol. It is stored at -20° C, wrapped in aluminum foil. Petri plates containing Tc are also wrapped in foil. Minimum inhibitory concentration is 3.1 ug/ml using a cell inoculum of  $5 \times 10^7$  to  $10^8$  per plate. More stringent selections employ up to 50 ug/ml Tc. When used for selection of plasmid maintenance, Tc concentration is 12.5 ug/ml.

b. Galactose utilization. Minimal A Medium (MILL72, p432), with galactose as carbon source: after autoclaving add (per liter) 1 ml 1 M MgSO<sub>4</sub>, 0.5 ml of 10 mg/ml thiamine HCl, 10 ml of 20% galactose, and amino acids as required. Cell inoculum per plate is less than  $5 \times 10^7$ .

c. Tc resistance and galactose utilization. Medium A with galactose (section b. above) is supplemented with Tc at 3.1 ug/ml.

#### Selectable systems for DBP isolation:

[0255] The tet gene from pBR322 and the *E. coli galT.K* genes are used in a gal deletion host strain for selection of DBP function. pKK175-6 (BROS84; Pharmacia, Piscataway, NJ), a pBR322 derivative, contains the replication origin, bla (confers Ap<sup>R</sup>) for selection of plasmid maintenance, and tet, one of the two selectable genes (Figure 3.) In pKK175-6, tet is promoterless, and all DNA upstream of the pBR322 tet coding region that potentially allow transcription in both directions (BROS82) have been deleted and replaced by the M13 mp8 polylinker. The polylinker and tet are flanked by strong transcription terminators from *E. coli rnaB*. tet is placed under control of the Tn5 neo promoter, P<sub>neo</sub>.

[0256] Plasmid pAA3H (figure 4) (ATCC #37,308) (AHME84) provides the second set of selectable genes, galT,K. In gal deleted hosts (such as strain ATCC #35,882 carrying the delta4 deletion (*E. coli delta4*)) plasmid pAA3H confers the Ap<sup>R</sup> Tc<sup>S</sup> Gal<sup>S</sup> phenotype (AHME84) because part of galE is deleted. The galT and galK genes in pAA3H are transcribed from the P<sub>1</sub> "antitet" promoter (BROS82). In *E. coli* strains carrying galT or galK mutations (e.g. strain HB101), pAA3H confers Gal<sup>+</sup>. We place galT<sup>+</sup> and galK<sup>+</sup> under control of the pBR322 amp gene promoter.

[0257] For both tet and gal systems, positive selections are used to select cells that either express or do not express these genes from cultures containing a vast excess of cells of the opposite phenotype.

#### Placement of test DNA binding sequence:

[0258] The test DNA binding sequence for the IDBP, λ O<sub>R</sub>3 (KIMJ87), is placed so that the first 5' base is the +1 base of the mRNA transcribed in each of the tet and gal transcription units (Table 100 and Table 101).

#### Engineering the idbp gene:

[0259] A DNA sequence encoding the wild-type Cro protein is designed such that expression is controlled by the lacUV5 promoter. The DNA sequence departs from the wild-type cro gene sequence by the introduction of restriction sites. Thus, the gene is called rav. The transcriptional unit comprising PlacUV5, rav, and trpA terminator is shown in Table 102.

#### Vector construction:

[0260] The construction of an operative cloning vector is summarized in Figure 5. The gal region of pAA3H requires manipulation before insertion into pKK175-6. First the λ-derived DNA between HpaI and EcoRI is replaced with a Clal linker (New England BioLabs, #1037). Standard methods are used and the resulting plasmid is named pEP1001 (Figure 6). All plasmids cited in the present application are catalogued in Table 103.

[0261] Next, we insert a synthetic fragment, shown in Table 104, comprising the phage fd terminator and two restriction sites (SpeI and SfiI) into the Clal site of pEP1001; the resulting plasmid is named pEP1002 (Figure 7). Next, we replace the P<sub>1</sub> promoter upstream of gal with P<sub>amp</sub> from pBR322. As shown in Table 100, λ O<sub>R</sub>3 is positioned downstream of P<sub>amp</sub> so that it can be used to determine whether binding of Cro can prevent transcription of galT,K. Restriction sites are provided to allow later alteration of the target sequence. The synthetic fragment is cloned into pEP1002 between DraIII and BamHI. The resulting plasmid is named pEP1003 and confers Gal<sup>S</sup> on delta4 cells.

[0262] The gal genes with the promoter and the fd terminator are moved from pEP1003 into pKK175-6. The 2.69 kb galT,K-bearing HpaI fragment of pEP1003 is ligated to DNA obtained from pKK175-6 by partial DraI digestion. Gal<sup>+</sup> colonies of transformed HB101 cells are picked. The resulting plasmid is named pEP1004 (Figure 9).

[0263] The Tn5 neo gene promoter and O<sub>R</sub>3 are synthesized (Table 101) and inserted upstream of the tet coding region of pEP1004 between the unique HindIII and SmaI sites. Plasmid DNA from Ap<sup>R</sup> Tc<sup>R</sup> Gal<sup>S</sup> colonies of transformed delta4 cells is analyzed for an insert in the EcoRI-EcoRV fragment of pEP1004. The resulting 7.1 kb plasmid, with two separate selectable gene systems under control of two different promoters and the test DNA binding sequence, is designated pEP1005 (Figure 10).

#### Cloning the idbp gene:

[0264] The BamHI site in the tet gene is removed from the tet gene in pEP1005 by site-directed mutagenesis; the sequence TGG-ATC-CTC that codes for W97-I98-L99 is changed to TGG-ATA-TTG. DNA from pEP1005 is linearized with EcoRV and part (ca. 10%) of the DNA is made single stranded with exonuclease III. The mutagenic oligonucleotide shown in Table 105 is annealed to the DNA that is then completed with Klenow enzyme and ligated. Plasmid DNA from Tc<sup>R</sup>, Gal<sup>+</sup> colonies of transformed HB101 is analyzed by standard means; the resulting plasmid is named pEP1006.

[0265] Synthetic DNA containing a SpeI overhang, followed by sequences for the lacUV5 promoter, a ribosome binding site, cloning sites for idbp, the trpA terminator (ROSE79), and an SfiI restricted end complementary to the SfiI site in pEP1006 is synthesized as six oligonucleotides as shown in Table 107. We use the methods of THER88 to anneal and ligate these fragments into SpeI, SfiI cut pEP1006. Plasmid DNA from Ap<sup>R</sup>, Tc<sup>R</sup>, Gal<sup>S</sup> colonies of transformed delta4 cells is examined for the SpeI-SfiI insertion by restriction with SpeI, BstEII, BglII, KpnI, and SfiI. The inserted DNA is verified by DNA sequencing, and the 7.22 kb plasmid containing the proper insertion is designated pEP1007, shown in Figure 11.

[0266] The idbp gene sequence specifying the Cro<sup>+</sup> protein and designated rav in this Example, is inserted in two cloning steps. The BstEII-BglII segment of rav (Table 109) is inserted first. Oligonucleotides olig#14 and olig#15 are synthesized, annealed, and filled in with Klenow enzyme (Cf. KARN84). The dsDNA is cut with BstEII and BglII and

ligated to BstEII-BglII cut pEP1007. The plasmid containing the appropriate partial ray sequence is designated pEP1008.

[0267] The BglII-KpnI fragment of ray is synthesized and inserted in the same manner as the BstEII-BglII fragment. (See Table 110.) This plasmid carrying the complete ray gene is designated pEP1009, shown in Figure 12.

Determine whether IDBP is expressed:

[0268] To determine whether cells carrying pEP1009 display the phenotypes expected for ray expression, the delta4 strain bearing pEP1009 is tested on various Ap containing selective media with and without IPTG. Cells are streaked on LB agar media containing: a) Tc; b) fusaric acid; or c) galactose (vide supra). Control strains are the delta4 host with no plasmid, and with pEP1005, pBR322, or pAA3H.

[0269] The results below indicate that the ray gene is expressed and the gene product is functional, and that expression is regulated by the lacUV5 promoter.

Growth of derivatives of strain delta4 on selective media (+ Ap)

[0270]

| supplements: | tetracycline |   | fusaric acid galactose |   |   |   |
|--------------|--------------|---|------------------------|---|---|---|
| IPTG:        | +            | - | +                      | - | + | - |
| plasmid:     |              |   |                        |   |   |   |
| -            | -            | - | -                      | - | - | - |
| pBR322       | +            | + | -                      | - | + | + |
| pAA3H        | -            | - | +                      | + | - | - |
| pEP1009      | -            | + | +                      | - | + | - |
| pEP1005      | +            | + | -                      | - | - | - |

[0271]  $\lambda$  cl<sup>-</sup> phage is streaked on each of the above strains, on LB agar with Ap, and with and without IPTG. At sufficiently high intracellular levels of Cro protein, binding of the Cro repressor protein to the  $\lambda$  phage operators  $O_R$  and  $O_L$  prevents phage growth. Data indicating correct expression and function of the ray gene are:

Growth of  $\lambda$  cl<sup>-</sup> on delta4 cells

[0272]

| plasmid | phage growth |       |
|---------|--------------|-------|
|         | +IPTG        | -IPTG |
| -       | +            | +     |
| pEP1009 | -            | +     |
| pEP1005 | +            | +     |

[0273] These procedures indicate that the chosen IDBP, the product of the ray gene, is expressed and is successfully repressing both the test operators on the plasmid and the wild type operators on the challenge phage.

DBP purification:

[0274] Proteins are purified as described by Leighton and Lu (LEIG87).

Quantitation of DBP binding:

[0275] We measure DBP binding to the target operator DNA sequence with a filter binding assay, initially using filter binding assay conditions similar to those described for  $\lambda$  Cro (KIMJ87). Data are analyzed by the methods of RIGG70 and KIMJ87.

[0276] The target DNA for the assay is the 113 bp ApaI-PsaI fragment from plasmid pEP1009 containing  $\lambda$  O<sub>R</sub>3. A control DNA fragment of the same size, used to determine non-specific DNA binding, contains a synthetic ApaI-XbaI DNA fragment specifying the amp promoter and the sequence

5' CTTATACACGAAGCGTGACAA 3'. This sequence preserves the base content of the O<sub>R</sub>3 sequence but lacks several sites of conserved sequence required for  $\lambda$  Cro binding (KIMJ87) and is cloned between the ApaI-XbaI sites of the pEP1009 backbone to yield pEP1010.

Media Formulations:

[0277] Gal<sup>S</sup> is demonstrable in LB agar and broth at very low concentrations (0.2% galactose), and is optimal at 2 to 8% galactose. Galactose and Tc selections are performed in LB medium. Fus<sup>S</sup> is best achieved in the medium described by Maloy and Nunn (MALO81) for E. coli K-12 strains.

Induction of DBP expression:

[0278] The pdbp gene is regulated by the lacUV5 promoter. Optimal induction is achieved by addition of IPTG at  $5 \times 10^{-4}$  M (MAUR80). Experimentation for each successful DBP determines the lowest concentration that is sufficient to maintain repression of the selection system genes.

Optimization of selections

[0279] For each selective medium used to detect IDBP function, factors are varied to obtain a maximal number of transformants per plate and with a minimal number of false positive artifactual colonies. Of greatest importance in this optimization is the transcriptional regulation of the initial potential-DBP, such that in further mutagenesis studies, de novo binding at an intermediate affinity is compensated by high level production of DBP.

Regulation of IDBP:

[0280] Cells carrying pEP1009 are grown in LB broth with IPTG at  $10^{-6}$ ,  $5 \times 10^{-6}$ ,  $10^{-5}$ ,  $5 \times 10^{-5}$ ,  $10^{-4}$  and  $5 \times 10^{-4}$  M. Samples are plated on LB agar and on LB agar containing fusaric acid or galactose as described in above. All media contain 200 ug/ml Ap, and the IPTG concentration of the broth culture media are maintained in the respective selective agar media.

[0281] The IPTG concentration at which 50% of the cells survive is a measure of affinity between IDBP and test operator, such that the lower the concentration, the greater the affinity. A requirement for low IPTG, e.g.  $10^{-6}$  M, for 50% survival due to Rav protein function suggests that use of a high level, e.g.  $5 \times 10^{-4}$  M IPTG, employed in selective media to isolate mutants displaying de novo binding of a DBP to target DNA, will enable isolation of successful DBPs even if the affinity is low.

Concentration of selective agents and cell inoculum size:

[0282] Fusaric acid and galactose content of each medium is varied, to allow the largest possible cell sample to be applied per Petri plate. This objective is obtained by applying samples of large numbers of sensitive cells (e.g.  $5 \times 10^7$ ,  $10^8$ ,  $5 \times 10^8$ ) to plates with elevated fusaric acid or galactose. Resistant cells are then used to determine the efficiency of plating. An acceptable efficiency is 80% viability for the resistant control strain bearing pEP1009 in a delta4 background. The total cell inoculum size is increased as is the level of inhibitory compound until viability is reduced to less than 80%.

Choice and cloning of target sequences:

[0283] Sequences of the human immunodeficiency virus type 1 (HIV-1) genome were searched for potential target sequences. The known sequences of isolates of HIV-1 were obtained from the GENBANK version 52.0 DNA sequence data base. First we found non-variable regions of HIV-1. We examined the HIV-1 genome from the TATA sequence in the 5'LTR of the HIV-1 genome to the end of the sequence coding for the tat and trg second exons. We intended to locate

non-variable regions where a DBP can interfere with the production of *tat* and/or *trc* mRNA because the products of these genes are essential in production of virus (DAYT86, FEIN86).

[0284] HIV-1 isolate HXB2 (RATN85) from nucleotide number 1 through 6100 is the reference to which we aligned all other HIV-1 isolates using the Nucleic Acid Database Search program (derived from FASTN (LIPM85)) in the IBI/Pustell Sequence Analysis Programs software package (International Biotechnologies, Inc., New Haven, CT). All stretches of at least 20 bases which have no variation in sequence among all HIV-1 isolates were retained as targets.

[0285] From the alignment, segments of the HIV-1 isolate HXB2 sequence that are non-variable among all HIV-1 sequences searched are:

|              |              |              |              |
|--------------|--------------|--------------|--------------|
| 350 - 371,   | 519 - 545,   | 623 - 651,   | 679 - 697,   |
| 759 - 781,   | 783 - 805,   | 1016 - 1051, | 1323 - 1342, |
| 1494 - 1519, | 1591 - 1612, | 1725 - 1751, | 1816 - 1837, |
| 2067 - 2094, | 2139 - 2164, | 2387 - 2427, | 2567 - 2606, |
| 2615 - 2650, | 2996 - 3018, | 3092 - 3117, | 3500 - 3523, |
| 3866 - 3887, | 4149 - 4170, | 4172 - 4206, | 4280 - 4302, |
| 4370 - 4404, | 4533 - 4561, | 4661 - 4695, | 4742 - 4767, |
| 4808 - 4828, | 4838 - 4864, | 4882 - 4911, | 4952 - 4983, |
| 5030 - 5074, | 5151 - 5173, | 5553 - 5573, | 5955 - 5991  |

[0286] In the present Example, these potential regions were searched for subsequences matching the central seven base pairs of the  $\lambda$  operators that have high affinity for  $\lambda$  Cro (*viz.*  $O_{R3}$ , the symmetric consensus, and the Kim *et al.* consensus (KIMJ87)). The consensus sequence of Kim *et al.* has higher affinity for Cro than does  $O_{R3}$  which is the natural  $\lambda$  operator having highest affinity for Cro. Cro is thought to recognize seventeen base pairs, with side groups on alpha 3 directly contacting the outer four or five bases on each end of the operator. Because the composition and sequence of the inner seven base pairs affect the position and flexibility of the outer five base pairs to either side, these bases affect the affinity of Cro for the operator.

[0287] The sequences sought are shown in Table 111. The letters "A" and "S" stand for antisense and sense. " $O_{R3A}$ /Symm. Consensus.5" is a composite that has  $O_{R3A}$  at all locations except 5, where it has the symmetric consensus base, C. Similarly, " $O_{R3A}$ /Symm. Consensus.6" has the symmetric consensus base at location 6 and  $O_{R3A}$  at other locations.

[0288] A FORTRAN program searched the non-variable HIV-1 subsequence segments for stretches of seven nucleotides of which at least five are G or C and which are flanked on either side by five bases of non-variable HIV-1 subsequence. The 427 candidate seven-base-pair subsequences obtained using these constraints on GC content were then searched for matches to either the sense or anti-sense strand sequences of the five seven-base-pair subsequences listed above. None of the HIV-1 subsequences is identical to any of the seven-base-pair subsequences. Three HIV-1 subsequences, shown in Table 112, were found that match six of seven bases. Eight subsequences, shown in Table 113, were found that match five out of seven bases and that have five or more GC base pairs. These HIV-1 subsequences are less preferred than the HIV-1 subsequences that match six out of seven bases.

5  
 12345678901234567  
 5' **a**ctTTccGCTggGgaCt **Bases 353-369**  
 actttccGCTggaaaagt **Left symmetrized**  
 10 agtccccGCTggggact **Right symmetrized**  
 tatcAcCGCAAgGgata **O<sub>R</sub>3**

15 (Lower case letters are palindromic in the two halves of the targets and O<sub>R</sub>3; highly conserved bases are bold and marked thus a.) Among the outer five bases of each half operator, bases 1 and 3 are palindromically related to bases 17 and 15 in Target HIV 353-369.

20  
 TCTCGAcGCAgGACTCG **Bases 681-697**  
 tctcgaCGCAgGcgaga **Left symmetrized**  
 cgagtAcGCAgGactcg **Right symmetrized**  
 25 tatcAcCGCAAgGgata **O<sub>R</sub>3**

None of bases 1-5 are palindromically related to bases 13-17 in Target HIV 681-697.

30  
 TTTGAcTAGCGgAGGCT **Bases 760-776**  
 tttgaCTAGCGgtcaaa **Left symmetrized**  
 35 agcctcTAGCGgaggct **Right symmetrized**  
 tatcAcCGCAAgGgata **O<sub>R</sub>3**

40 None of bases 1-5 are palindromically related to bases 13-17 in Target HIV 760-776.

[0289] There is extensive sequence variability among the twelve phage  $\lambda$  operator half-sites. For example:

45  
 tAtCaCCGCCGGtGaTa **Consensus**  
 tAtCaCCGCCaAgGaTa **O<sub>R</sub>3A**

50 The bases in lower case in Consensus and O<sub>R</sub>3 sequences shown above are more variable among various lambdoid operators than are bases shown by upper case letters. Studies of mutant operators indicate that A2 and C4 are required for Cro binding. In Target HIV 353-369, bases T3, C6, C7, G8, C9, G14, and A15 match the symmetric consensus sequence, but the highly conserved A2 and C4 are different from lambdoid operators and Cro will not bind to these sub-sequences. Mutagenesis of the DNA-contacting residues of alpha 3 is thus the first step in producing a DBP that recognizes the left symmetrized or right symmetrized target sequences.

[0290] Target HIV 353-369 is a preferred target because the core (underlined above) is highly similar to the Kim et al. consensus. Target HIV 760-776 is preferred over Target HIV 681-697 because it is highly similar to O<sub>R</sub>3.

[0291] The method of the present invention does not require any similarity between the target subsequence and



the original binding site of the initial DBP. The fortuitous existence of one or more subsequences within the target genes that has similarity to the original binding site of the initial DBP reduces the number of iterative steps needed to obtain a protein having high affinity and specificity for binding to a site in the target gene.

[0292] Since the target sequence is from a pathogenic organism, we require that the chosen target subsequence be absent or rare in the genome of the host organism, e.g. the target subsequences chosen from HIV should be absent or rare in the human genome.

[0293] Candidate target binding sites are initially screened for their frequency in primate genomes by searching all DNA sequences in the GENBANK Primate directory (2,258,436 nucleotides) using the IBI/Pustell Nucleic Acid Database Search program to locate exact or close matches. A similar search is made of the E. coli sequences in the GENBANK Bacterial directory and in the sequence of the plasmid containing the idbp gene. The sequences of potential sites for which no matches are found are used to make oligonucleotide probes for Southern analysis of human genomic DNA (SOUT75). Sequences which do not specifically bind human DNA are retained as target binding sequences.

[0294] The HIV 353-369 left symmetrized and right symmetrized target subsequences are inserted upstream of the selectable genes in the plasmid pEP1009, replacing the test sequences, to produce two operative cloning vectors, pEP1011 and pEP1012, for development of Rav<sub>L</sub> and Rav<sub>R</sub> DBPs. The promoter-test sequence cassettes upstream of the tet and gal operon genes are excised using StuI-HindIII and ApaI-XbaI restrictions, respectively. Replacement promoter-target sequence cassettes are synthesized and inserted into the vector, replacing O<sub>R</sub>3 with the HIV 353-369 left or right symmetrized target sequence in the sequences shown in Table 100 and Table 101.

#### Choice of residues in Cro to vary:

[0295] The choice of the principal and secondary sets of residues depends on the goal of the mutagenesis. In the protocol described here we vary, in separate procedures, the residues: a) involved in DNA recognition by the protein, and b) involved in dimerization of the protein. In this section we identify principal and secondary sets of residues for DNA recognition and dimerization.

#### Pick principal set for DNA-recognition:

[0296] The principal set of residues involved in DNA-recognition is defined as those residues which contact the operator DNA in the sequence-specific DNA-protein complex. Although no crystal structure of a  $\lambda$  Cro-operator DNA complex is available, a crystal structure of a complex between the structural homolog 434 repressor N-terminal domain and a consensus operator has been described (ANDE87). A crystal structure of Cro dimer has been determined (ANDE81) and modeling studies have suggested residues that can make sequence-specific or sequence-independent contacts with DNA in sequence-specific complexes (TAKE83, OHLE83, TAKE85, TAKE86). Isolation and characterization of Cro mutants have identified residues which contact DNA in protein-operator complexes (PAKU86, HOCH86a,b, EISE85).

[0297] Important contacts with DNA are made by protein residues in and around the H-T-H region and in the C-terminal region. Hochschild et al. (HOCH86a,b) have presented direct evidence that Cro alpha helix 3 residues S28, N31, and K32 make sequence-specific contacts with operator bases in the major groove. Mutagenesis experiments (EISE85, PAKU86) and modeling studies (TAKE85) have implicated these residues as well. In addition, these studies suggest that H-T-H region residues Q16, K21, Y26, Q27, H35, A36, R38, and K39 also make contacts with operator DNA. In the C-terminal region, mutagenesis experiments (PAKU86) and chemical modification studies (TAKE86) have identified K56, and K62 as making contacts to DNA. In addition, computer modeling suggests that the 5 to 6 C-terminal amino acids of  $\lambda$  Cro can contact the DNA along the minor groove (TAKE85). From these considerations, we select the following set of residues as a principal set for use in variegation steps intended to modify DNA recognition by Cro or mutant derivative proteins: 16, 21, 26, 27, 31, 32, 35, 36, 38, 39, 56, 62, 63, 64, 65, 66.

#### Pick secondary set for DNA recognition:

[0298] The residues in the secondary set contact or otherwise influence residues in the principal set. A secondary set for DNA recognition includes the buried residues of alpha helix 3: A29, I30, A33, and I34. Interactions between buried residues in alpha helix 2 and buried residues in alpha helix 3 are known to stabilize H-T-H structure and residues in the turn between alpha helix 2 and alpha helix 3 of H-T-H proteins are conserved among these proteins (PTAS86 p102). In  $\lambda$  Cro these positions are T17, T19, A20, L23, G24, and V25. Changes in the dimerization region can influence binding. In  $\lambda$  Cro, residues thought to be involved in dimer stabilization are E54, V55, and F58 (TAKE85, PABO84). Finally, residues influencing the position of the C-terminal arm of  $\lambda$  Cro are P57, P59, and S60. Thus the secondary set of residues for use in variegation steps intended to modify DNA recognition by  $\lambda$  Cro or Rav proteins is: 17, 19, 20, 23, 24, 29, 30, 33, 34, 54, 55, 57, 58, 59 and 60.

Pick principal set for dimerization:

[0299] Different principal and secondary sets of residues must be picked for use in variegation steps intended to alter dimer interactions. In  $\lambda$  Cro, antiparallel interactions between E54, V55, and K56 on each monomer have been proposed to stabilize the dimer (PABO84). In addition, F58 from one monomer has been suggested to contact residues in the hydrophobic core of the second monomer. Inspection of the 3D structure of  $\lambda$  Cro suggests important contacts are made between F58 of one monomer and I40, A33, L23, V25, E54, and A52. In addition, residues L7, I30, and L42 of one monomer could make contact with a large side chain positioned at 58 in the other monomer. Thus, a set of principal residues includes: 7, 23, 25, 30, 33, 40, 42, 52, 54, 55, 56, and 58.

Pick secondary set for dimerization:

[0300] The secondary set of residues for variegation steps used to alter dimer interactions includes residues in or near the antiparallel beta sheet that contains the dimer forming residues. Residues in this region are E53, P57, and P59. Residues in alpha helix 1 influencing the orientation of principal set residues are K8, A11, and M12. Residues in the antiparallel beta sheet formed by the beta strands 1, 2, and 3 (see Table 1) in each monomer also influence residues in the principal set. These residues include I5, T6, K39, F41, V50, and Y51. Thus the set of secondary residues includes: 5, 6, 8, 11, 12, 41, 50, 51, 53, 57, and 59.

Pick the range of variation for alteration of DNA binding:

[0301] For the initial variegation step to produce a modified Rav protein with altered DNA specificity a set of 5 residues from the principal set is picked. Focused Mutagenesis is used to vary all five residues through all twenty amino acids. The residues are picked from the same interaction set so that as many as  $3.2 \times 10^7$  different DNA binding surfaces will be produced.

[0302] A number of studies have shown that the residues in the N-terminal half of the recognition helix of an H-T-H protein strongly influence the sequence specificity and strength of protein binding to DNA (HOCH86a,b, WHAR85, PABO84). For this reason we choose residues Y26, Q27, S28, N31, and K32 from the principal set as residues to vary in the first variegation step. Using the optimized nucleotide distribution for Focused Mutagenesis described above, and assuming that  $S_{err} = 5\%$  as defined at the start of this Example, the parental sequence is present in the variegated mixture at one part in  $3.1 \times 10^5$  and the least favored sequence, F at each residue, is present at one part in  $10^8$ . Thus, this level of variegation is well within bounds for a synthesis, ligation, transformation, and selection system capable of examining  $5 \times 10^8$  DNA sequences.

Pick the range of variation of residues for alteration of dimerization:

[0303] As described in the Detailed Description and in this Example, altered  $\lambda$  Cro proteins, Rav<sub>L</sub> and Rav<sub>R</sub>, that bind specifically and tightly to left and right symmetrized targets derived from HIV 353-369, are first developed through one or more variegation steps. Site-specific changes are then engineered into Rav<sub>L</sub> to produce dimerization defective proteins. Structure-directed Mutagenesis is performed on Rav<sub>R</sub> to produce mutations in Rav<sub>R</sub> that can complement dimerization defective Rav<sub>L</sub> proteins and produce obligate heterodimers that bind to HIV 353-369.

[0304] One of the interactions in the dimerization region of  $\lambda$  Cro is the hydrophobic contact between residues V55 of both monomers. The VF55 mutation substitutes a bulky hydrophobic side group in place of the smaller hydrophobic residue; other substitutions at residue 55 can be made and tested for their ability to dimerize. A small hydrophobic or neutral residue present at residue 55 in a protein encoded on expression by a second gene may result in obligate complementation of VF55. In addition, changes in nearby components of the beta strand, E53, E54, K56, and P57 may effect complementation. Thus a set of residues for the initial variegation step to alter the Rav<sub>R</sub> dimer recognition is 53, 54, 55, 56, and 57.

[0305] Another interaction in the dimerization region of  $\lambda$  Cro is the hydrophobic contact between F58 of one monomer with the hydrophobic core of the other monomer. As mentioned above residues L7, L23, V25, A33, I40, L42, A52, and E54 of one monomer all could make contacts with a large residue at position 58 in the other monomer. The FW58 mutation inserts the largest aromatic amino acid at this position. Compensation for this substitution may require several changes in the hydrophobic core of the complementing monomer. Residues for Focused Mutagenesis in the initial variegation step to alter Rav<sub>R</sub> dimer recognition in this case are: 23, 25, 33, 40, and 42.

[0306] In each of the two cases described above, the initial variegation step involves Focused Mutagenesis to alter 5 residues through all twenty amino acids. As was shown in Section 6.2.5, this level of variegation is within the limits set by using optimized codon distributions and the values for  $S_{err}$  and transformation yield assumed at the start of this Example.

Mutagenesis of DNA:

[0307] Codons encoding  $\lambda$  Cro residues Y26, Q27, S28, N31, and K32 are contained in a 51 bp PpuMI to BglII fragment of the rav gene. To produce the cassette containing the variegated codons we synthesize the 66 nucleotide anti-sense variegated strand, olig#50, and the primer, olig#52:

```

      d   l   g   v   X   X   X   a   i   X
      22  23  24  25  26  27  28  29  30  31
10  5' t cct aAG GAC CTA GGG GTG fzk fzk fzk GCG ATT fzk
      | PpuMI |
15
      X   a   i   h   a   g   r   k   i
      32  33  34  35  36  37  38  39  40
20  fzk GCC ATC CAT GCC GGC CGA AAG ATC Tt 3' olig#50
      3'-ccg gct ttc tag aacgccgtg-5' olig#52
      | BglII |

```

The position of the amino acid residue in  $\lambda$  Cro is shown above the codon for the residue. Unaltered residues are indicated by their lower case single letter amino acid codes shown above the position number. Variegated residues are denoted with an upper case, bold X. The restriction sites for PpuMI, and BglII are indicated below the sequence. Since restriction enzymes do not cut well at the ends of DNA fragments, 5 extra nucleotides have been added to the 5' end of the cassette. These extra nucleotides are shown in lower case letters and are removed prior to ligating the cassette into the operative vector. The sequence "fzk" denotes the variegated codons and indicates that nucleotide mixtures optimized for codon positions 1, 2, or 3 are to be used. "t" is a mixture of 26% T, 18% C, 26% A, and 30% G, producing four possibilities. "z" is a mixture of 22% T, 16% C, 40% A, and 22% G, producing four possibilities. "k" is an equimolar mixture of T and G, producing two possibilities. Each "fzk" codon produces  $4 \times 4 \times 2 = 2^5 = 32$  possible DNA sequences, coding on expression for 20 possible amino acids and stop. The DNA segment above comprises  $(2^5)^5 = 2^{25} = 3.2 \times 10^7$  different DNA sequences coding on expression for  $20^5 = 3.2 \times 10^6$  different protein sequences.

[0308] After synthesis and purification of the variegated DNA, the oligonucleotides #50 and #52 are annealed and the resulting superoverhang is filled in using Klenow fragment as described by Hill (AUSU87, Unit 8.2). The double stranded oligonucleotide is digested with the enzymes PpuMI and BglII and the mutagenic cassette is purified as described by Hill. The mutagenic cassette is cloned into the vectors pEP1011 and pEP1012 which have been digested with PpuMI and BamHI, and the ligation mixtures containing variegated DNA are used to transform competent delta4 cells. The transformed cells are selected for vector uptake and for successful repression at low stringency as described above. Cells containing Rav proteins that bind to the left or right symmetrized targets display the  $Tc^S$ ,  $Fus^R$  and  $Gal^R$  phenotypes.

[0309] Surviving colonies are screened for correct  $DBP^+$  and  $DBP^-$  phenotypes in the presence or absence of IPTG as described above. Relative measures of the strengths of  $DBP$ -DNA interactions in vivo are obtained by comparing phenotypes exhibited at reduced levels of IPTG.  $DBP$  genes from clones exhibiting the desirable phenotypes are sequenced. Plasmid numbers from pEP1100 to pEP1199 are reserved for plasmids yielding rav<sub>L</sub> genes encoding proteins that bind to the Left Symmetrized Targets carried on the plasmids. Similarly, plasmid numbers pEP1200 through pEP1299 plasmids containing rav<sub>R</sub> genes encoding proteins that bind to the Right Symmetrized Targets carried on these plasmids.

[0310] Based on the determinations above, one or more  $Rav_L$  and  $Rav_R$  proteins are chosen for further analysis in vitro. Proteins are purified as described above. Purified  $DBPs$  are quantitated and characterized by absorption spectroscopy and polyacrylamide gel electrophoresis.

[0311] In vitro measurements of protein-DNA binding using purified  $DBPs$  are performed as described in the Overview: DNA-Binding, Protein Purification, and Characterization and in this Example. These measurements determine equilibrium binding constants ( $K_D$ ), and the dissociation ( $k_d$ ) and association ( $k_a$ ) rate constants for sequence-specific and sequence-independent  $DBP$ -DNA complexes. In addition, DNase protection assays are used to demonstrate spe-

cific DBP binding to the Target sequences.

[0312] Estimates of relative DBP stability are obtained from measurements of the thermal denaturation properties of the proteins. *In vitro* measures of protein thermal stability are obtained from determinations of protein circular dichroism and resistance to proteolysis by thermolysin at various temperatures (HECH84) or by differential scanning calorimetry (HECH85b).

[0313] One or more iterations of variegation, involving residues thought capable of influencing DNA binding, of the *rav<sub>L</sub>* and *rav<sub>R</sub>* genes produce *Rav<sub>L</sub>* and *Rav<sub>R</sub>* proteins that bind tightly and specifically to the HIV 353-369 left and right symmetrized targets. Additional variegation steps, to optimize protein binding properties can be performed as outlined in the Overview: Variegation Strategy.

[0314] By hypothesis, we isolate pEP1127 that contains a *pdbp* gene that codes on expression for *Rav<sub>L</sub>*-27, shown in Table 114, that binds the left-symmetrized target best among selected *Rav<sub>L</sub>* proteins. Similarly, pEP1238 contains a *pdbp* gene that codes on expression for *Rav<sub>R</sub>*-38, shown in Table 115, that binds the right-symmetrized target best among selected *Rav<sub>R</sub>* proteins.

[0315] We now use the genes for the *Rav<sub>R</sub>* and *Rav<sub>L</sub>* monomers as starting points for production of obligately heterodimeric proteins *Rav<sub>L</sub>*:*Rav<sub>R</sub>* that recognize the HIV 353-369 target. First we change the target sequences in pEP1238 (containing *rav<sub>R</sub>*-38). We replace both occurrences of the Right Symmetrized Target (in *tet* and *galT<sub>K</sub>* promoters) with the HIV 353-369 target sequence. *Delta4* cells containing plasmids carrying the HIV 353-369 targets display the *Ap<sup>R</sup>*, *Tc<sup>R</sup>*, *Fus<sup>S</sup>* and *Gal<sup>S</sup>* phenotypes. Plasmids carrying HIV 353-369 targets and the *rav<sub>R</sub>* gene are designated by numbers pEP1400 through pEP1499 and corresponding to the number of the donor plasmid of the 1200 series; for example, replacing the target sequences in pEP1238 produces pEP1438.

#### Engineering dimerization mutants of *Rav<sub>L</sub>*:

[0316] To create the site specific VF55 and FW58 mutations in *rav<sub>L</sub>* we synthesize the two mutagenesis primers:

|        |     |     |     |            |          |     |     |    |      |
|--------|-----|-----|-----|------------|----------|-----|-----|----|------|
|        |     | a   | e   | e          | f        | k   | p   | f  |      |
|        |     | 52  | 53  | 54         | 55       | 56  | 57  | 58 |      |
| 5'     | GGC | GAA | GAG | <u>TTC</u> | AAG      | CCC | TTC | 3' | VF55 |
| primer |     |     |     |            |          |     |     |    |      |
|        |     |     |     |            |          |     |     |    |      |
|        |     | v   | k   | p          | <u>W</u> | p   | s   | n  |      |
|        |     | 55  | 56  | 57         | 58       | 59  | 60  | 61 |      |
| 5'     | GTA | AAG | CCC | <u>TGG</u> | CCC      | AGT | AAC | 3' | FW58 |
| primer |     |     |     |            |          |     |     |    |      |

Underlining indicates the varied codons and residues. The plasmid pEP1127 (containing *rav<sub>L</sub>*-27) is chosen for mutagenesis. The gene fragment coding on expression for the carboxy-terminal region of the *Rav<sub>L</sub>* protein is transferred into M13mp18 as a *Bam*HI to *Kpn*I fragment. Oligonucleotide-directed mutagenesis is performed as described by Kunkel (AUSU87, Unit 8.1). The fragment bearing the modified region of *Rav<sub>L</sub>* is removed from M13 RF DNA as the *Bam*HI to *Kpn*I fragment and ligated into the correct location in the pEP1100 vector. Mutant-bearing plasmids are used to transform competent cells. Transformed cells are selected for plasmid uptake and screened for DBP<sup>+</sup> phenotypes (*Tc<sup>R</sup>*, *Fus<sup>S</sup>*, and *Gal<sup>S</sup>* in *E. coli delta4*; *Gal<sup>+</sup>* in *E. coli* HB101). Plasmids isolated from DBP<sup>+</sup> cells are screened by restriction analysis for the presence of the *rav<sub>L</sub>* gene and the site-specific mutation is confirmed by sequencing. The plasmid containing the *rav<sub>L</sub>*-27 gene with the VF55 mutation is designated pEP1301. Plasmid pEP1302 contains the *rav<sub>L</sub>*-27 gene with the FW58 alteration.

[0317] For the production of obligate heterodimers as described below, the *rav<sub>L</sub>*- genes encoding the VF55 or FW58 mutations are excised from pEP1301 or pEP1302 and are transferred into plasmids containing the gene for Km and neomycin resistance (*neo*, also known as *not I*). These constructions are performed in three steps as outlined below. First, the *neo* gene from Tn5 coding for Km<sup>R</sup> and contained on a 1.3 Kbp *Hind*III to *Sma*I DNA fragment is ligated into the plasmid pSP64 (Promega, Madison, WI) which has been digested with both *Hind*III and *Sma*I. The resulting 4.3 kbp plasmid, pEP1303, confers both Ap and Km resistance on host cells. Next, the *bla* gene is removed from pEP1303 by digesting the plasmid with *Aat*II and *Bgl*II. The 3.5 Kbp fragment resulting from this digest is purified, the 3' overhang-

ing ends are blunted using T4 DNA polymerase (AUSU87, Unit 3.5), and the fragment is recircularized. This plasmid is designated pEP1304 and transforms cells to Km resistance. In the final step, the *rav<sub>L</sub>* gene is incorporated in to pEP1304. Plasmid pEP1301 or pEP1302 is digested with *Sfi*I and the resulting 3' overhangs are blunted using T4 DNA polymerase. Next the linearized plasmid is digested with *Spe*I and the resulting 5' overhangs are blunted using the Klenow enzyme reaction (KLEN70). The ca. 340 bp blunt-ended DNA fragment containing the entire *rav<sub>L</sub>* gene is purified and ligated into the *Pvu*II site in pEP1304. Transformed cells are selected for Km<sup>R</sup> and screened by restriction digest analysis for the presence of *rav<sub>L</sub>* genes. The presence of *rav<sub>L</sub>* genes containing the site-specific VF55 or FW58 mutations is confirmed by sequencing. The plasmid containing the *rav<sub>L</sub>* gene with the VF55 mutation is designated pEP1305. The plasmid containing the *rav<sub>L</sub>* gene with the FW58 mutations is designated pEP1306.

[0318] In a manner similar to the constructions described above, we ligate the original unmodified *rav<sub>L</sub>* gene into pEP1304 to produce plasmid pEP1307.

#### Engineering heterodimer binding of target DNA:

[0319] This round of variegation is performed to produce mutations in *Rav<sub>R</sub>* proteins that complement the dimerization deficient mutations in the *Rav<sub>L</sub>* proteins produced above. To complement the FW58 mutation, the set of five residues L23, V25, A33, I40, and L42 are chosen from the primary set of residues as targets for Focused Mutagenesis.

[0320] In an initial series of procedures to test for recognition of HIV 353-369 by the heterodimer *Rav<sub>L</sub>:Rav<sub>R</sub>*, we transform cells containing pEP1438 (containing *rav<sub>R</sub>*-38 and HIV 353-369 targets) with pEP1307 (containing *rav<sub>L</sub>*).

Intracellular expression of *rav<sub>L</sub>* and *rav<sub>R</sub>* produces a population of dimeric repressors: *Rav<sub>L</sub>:Rav<sub>L</sub>*, *Rav<sub>L</sub>:Rav<sub>R</sub>* and *Rav<sub>R</sub>:Rav<sub>R</sub>*. If the heterodimeric protein is formed and binds to HIV 353-369, cells expressing both *rav* alleles will exhibit the Km<sup>R</sup> Ap<sup>R</sup> Gal<sup>R</sup> Fus<sup>R</sup> phenotypes (*vide infra*). Several pairs of *rav<sub>L</sub>* and *rav<sub>R</sub>* genes are used in parallel procedures; the best pair is picked for use and further study. Selections for binding the HIV 353-369 target by the heterodimeric protein can be optimized using this system.

[0321] Focused Mutagenesis of residues 23, 25, 33, 40, and 42 requires the synthesis and annealing of two overlapping variegated strands because in the *rav* gene a single cassette spanning these residues extends from the *Bal*I site to the *Bam*HI site and exceeds the assumed synthesis limit of 100 nucleotides. As no variegation affects the overlap, the annealing region is complementary. The antisense strand of the DNA sequence from the *Bal*I site blunt end to the end of the codon for G37 is denoted olig#53.

```

      q   t   k   t   a   k   d   X   g   X   y   q
      16  17  18  19  20  21  22  23  24  25  26  27
5' C CAA ACC AAG ACA GCG AAG GAC fzk GGG fzk TAT CAG
      |BalI|
      s   a   i   n   k   X   i   h   a   g
      28  29  30  31  32  33  34  35  36  37
AGC GCG ATT AAC AAG fzk ATC CAT GCC GGC 3' olig#53

```

f = (26% T, 18% C, 26% A, 30% G)

z = (22% T, 16% C, 40% A, 22% G)

k = equimolar T and G

Olig#53 contains vg codons for residues 23, 25, and 33.

[0322] Olig#54 is the sense strand from base 1 in codon 34 to the *Bam*HI site:

i h a g r k X  
 34 35 36 37 38 39 40  
 5 3' TAG GTA CGG CCG GCA TTC jqm  
  
 f X t i n a d n k  
 10 41 42 43 44 45 46 47 48 49  
 AAG jqm TGG TAA TTG CGA CTA CCT AGG cca ca 5' olig#54  
BamHI

15

j = (26% A, 18% G, 16% T, 30% C)  
 q = (22% A, 16% G, 40% T, 22% C)  
 20 m = equimolar A and C

25

Olig#54 contains variegated codons for residues 40 and 42. Since olig#54 is the sense strand, the variegated nucleotide distributions must complement the distributions for codon positions 1, 2, and 3 used in the antisense strand. These sense codon distributions are designated "j", "q", and "m", and represent the complements to the optimized codon distributions developed for codon positions 1, 2, and 3, respectively, in the antisense strand. The two strands (olig#53 and olig#54) share a 12 nucleotide overlap extending from the first position in the codon for I34 to the end of the codon for G37. The overlap region is 66% G or C.

30

[0323] The two strands shown above are synthesized, purified, annealed, and extended to form dsDNA. Following restriction endonuclease digestion and purification, the mutagenic cassettes are ligated into pEP1438 (containing the asymmetric HIV 353-369 target) in the appropriate locus in the *ray<sub>R</sub>* gene. The ligation mixtures are used to transform competent cells that contain pEP1306 (the plasmid with the *ray<sub>L</sub>* gene carrying the FW58 site-specific mutation).

35

[0324] Above we picked a set of five residues in  $\lambda$  Cro, E53, E55, V55, K56, and P57, as targets for focused mutagenesis in the first variegation step of the procedure to produce a *Ray<sub>R</sub>* protein that complements the dimerization-deficient VF55 *Ray<sub>L</sub>* mutation. These five residues are contained on a 71 bp BamHI to KpnI fragment of the *ray* gene (Table 100). To produce a cassette containing the variegated codons we synthesize olig#58:

40

g s v y a X X X X X f  
 48 49 50 51 52 53 54 55 56 57 58  
 5' ct gat GGA TCC GTC TAC GCG fzk fzk fzk fzk fzk TTC  
BamHI

45

p s n k k  
 59 60 61 62 63  
 CCG AGT AAC AAA AAA

50

t t a .  
 64 65 66 67  
 ACA ACA GCG TAA TAGTAGGTACC ta 3' olig#58

55

KpnI

[0325] After synthesis and purification of the vgDNA, strands are self-annealed using the 10 nucleotide palindrome at the 3' end of the sequence. The resulting superoverhangs are filled in using the Klenow enzyme reaction as described previously and the double-stranded oligonucleotide is digested with *Bam*HI and *Kpn*I. Purified mutagenic cassettes are ligated into one or more operative vectors (picked from the pEP1200 series) in the appropriate locus in the *rav<sub>R</sub>* gene. The ligation mixtures are used to transform competent cells that contain pEP1305 (the plasmid carrying the *rav<sub>L</sub>* gene with the FV55 mutation).

[0326] Operative vectors carrying the VF55 or FW58 mutation in *rav<sub>L</sub>* confer Km resistance. Operative vectors carrying mutagenized *rav<sub>R</sub>* genes contain the gene for Ap<sup>R</sup> as well as the selective gene systems for the DBP<sup>+</sup> phenotypes. Cells containing complementing mutant proteins are selected by requiring both Ap<sup>R</sup> and Km<sup>R</sup> and repression of the complete HIV 343-369 target sequence (substituted for the Left and Right Symmetrized Targets in the selection genes). Cells possessing the desired phenotype are Ap<sup>R</sup>, Km<sup>R</sup>, Fus<sup>R</sup>, and Gal<sup>R</sup> (in *E. coli* delta4).

[0327] Plasmids from candidate colonies are first isolated genetically by transformation of cells at low plasmid concentration. Cells carrying plasmids coding for Rav<sub>L</sub> proteins will be Km<sup>R</sup>, while cells carrying plasmids coding for Rav<sub>R</sub> proteins will be Ap<sup>R</sup>. Plasmids are individually screened to ensure that they confer the DBP<sup>+</sup> phenotype and are characterized by restriction digest analysis to confirm the presence of *rav<sub>L</sub>* or *rav<sub>R</sub>* genes. Plasmid pairs are co-tested for complementation by restoration of the DBP<sup>+</sup> phenotype when both *rav<sub>R</sub>* and *rav<sub>L</sub>* are present intracellularly. Successfully complementing plasmids are sequenced through the *rav* genes to identify the mutations and to suggest potential locations for optional subsequent rounds of variegation.

[0328] Plasmids containing genes for altered Rav<sub>R</sub> proteins that successfully complement the *rav<sub>L</sub>* VF55 mutation are designated by plasmid numbers pEP1500 to pEP1599. Similarly, plasmids containing genes for altered Rav<sub>R</sub> proteins that successfully complement the *rav<sub>L</sub>* FW58 mutation are designated by plasmid numbers pEP1600 to pEP1699.

[0329] Heterodimeric proteins are purified and their DNA-binding and thermal stability properties are characterized as described above. Pairwise variation of the Rav<sub>R</sub> and Rav<sub>L</sub> monomers can produce dimeric proteins having different dimerization or dimer-DNA interaction energies. In addition, further rounds of variegation of either or both monomers to optimize DNA binding by the heterodimer, dimerization strength or both may be performed.

[0330] In this manner a heterodimeric protein that recognizes any predetermined target DNA sequence is constructed. The foregoing is hypothetical. The sequences shown as the result of selection are given by way of example and must not be construed as predictions that proteins of the stated sequence will have specific affinity for any DNA sequence.

## Example 2

[0331] Presented below is a hypothetical example of a protocol for developing new DNA-binding polypeptides, derived from the first ten residues of phage P22 Arc and a segment of variegated polypeptide with affinity for DNA subsequences found in HIV-1 using *E. coli* K12 as the host cell line. Some further optimization, in accordance with the teachings herein, may be necessary to obtain the desired results. Possible modifications in the preferred method are discussed immediately following the hypothetical example.

[0332] We set the same hypothetical technical capabilities as used in Detailed Example 1.

## Overview:

[0333] To obtain significant binding between a genetically encoded polypeptide and a predetermined DNA subsequence, the surfaces must be complementary over a large area, 1000 Å<sup>2</sup> to 3000 Å<sup>2</sup>. For the binding to be sequence-specific, the contacts must be spread over many (12 to 20) bases. An extended polypeptide chain that touches 15 base pairs comprises at least 25 amino acids. Some of these residues will have their side groups directed away from the DNA so that many different amino acids will be allowed at such residues, while other residues will be involved in direct DNA contacts and will be strongly constrained. Unless we have 3D structural data on the binding of an initial polypeptide to a test DNA subsequence, we can not *a priori* predict which residues will have their side groups directed toward the DNA and which will have their side groups directed outward. We also can not predict which amino acids should be used to specifically bind particular base pairs. Current technology allows production of 10<sup>7</sup> to 10<sup>8</sup> independent transformants per µg of DNA which allows variation of 5 or 6 residues through all twenty amino acids. Alternatively, between 23 and 30 two-way variations of DNA bases can be applied that will affect between 8 and 30 codons.

[0334] Sauer and colleagues (VERS87b) have shown that P22 Arc binds to DNA using a motif other than H-T-H. There is as yet no published X-ray structure of Arc, though the protein has been crystallized and diffraction data have been collected (JORD85). A combination of genetics and biochemistry indicates that the first 10 residues of each Arc monomer (M-K-G-M-S-K-M-P-Q-F) bind to palindromically related sets of bases on either side of the center of symmetry of the 21 bp operator shown in Table 200. Furthermore, the first ten residues of each Arc monomer assume an extended conformation (VERS87b). The hydrophobic residues may be involved in contacts to the rest of the protein, but

there are several examples from H-T-H DBPs of hydrophobic side groups being in direct contact with bases in the major groove. We do know that these first ten residues of Arc can exist in a conformation that makes sequence-specific favorable contacts with the arc operator.

[0335] We pick a target DNA subsequence from the HIV-1 genome such that a portion of the chosen sequence is similar to one half-site of the arc operator. We use part of this chosen sequence for an initial chimeric target. One half of the first target is the DNA subsequence obtained from HIV-1 and the other half of the target is one half-site of the arc operator. For this example, we will use a plasmid bearing wild-type arc operators repressed by the Arc repressor as a control. After demonstrating that Arc repressor can regulate the selectable genes, we replace the wild-type arc operator with the target DNA subsequence. We then replace the arc gene with a variegated pdcbp gene and select for cells expressing DBPs that can repress the selectable genes.

[0336] Once a protein is obtained that binds to the target that has similarity to one half of the Arc operator, we can change the target so that it has less similarity to one half of the Arc operator and mutagenize those residues that correspond to residues 1-10 of Arc. *In vivo* selection will isolate a protein that binds to the new target. A few repetitions of this process can produce a polypeptide that binds to any predetermined DNA sequence.

[0337] Our potential DNA-binding polypeptide (DBP) will be 36 residues long and will contain the first ten residues of Arc which are thought to bind to part of the half operator. DNA encoding the first ten amino acids of Arc is linked at the 3' terminus of this gene fragment to vgDNA that encodes a further 26 amino acids. Twenty-four of the codons encode two alternative amino acids so that  $2^{24}$  = approx.  $1.6 \times 10^7$  protein sequences result. The amino acids encoded are chosen to enhance the probability that the resulting polypeptide will adopt an extended structure and that it can make appropriate contacts with DNA. The Chou-Fasman (CHOU78a, CHOU78b) probabilities are used to pick amino acids with high probability of forming beta structures (M, V, I, C, F, Y, Q, W, R, T); the amino acids are grouped into five classes in Table 16. In addition, to discourage sequence-independent DNA binding, some acidic residues should be included. Glutamic acid is a strong alpha helix former, so in early stages we use D exclusively. Further, S and T both can make hydrogen bonds with their hydroxyl groups, but T favors extended structures while S favors helices; hence we use only T in the initial phase. Likewise, N and Q provide similar functionalities on their side groups, but Q favors beta and so is used exclusively in initial phases. Positive charge is provided by K and R, but only R is used in the variegated portion. Alanine favors helices and is excluded. P kinks the chain and is allowed only near the carboxy terminus in initial iterations.

[0338] After one selection, we design a different set of binary variegations that includes the selected sequence and perform a second mutagenesis and selection. After two or more rounds of diffuse variegation and selection, we choose a subset of residues and vary them through a larger set of amino acids. We continue until we obtain sufficient affinity and specificity for the target. None of the polypeptides discussed in this example is likely to have a defined 3D structure of its own, because they are all too short. Even if one folded into a definite structure, that structure is unlikely to be related to DNA-binding. A 3D structure, obtained by X-ray diffraction or NMR, of a DNA-polypeptide complex would give us useful indications of which residues to vary. Scattering the variegation along the chain and sampling different charges, sizes, and hydrophobicities produces a series of proteins, isolated by *in vivo* selection, with progressively higher affinity for the target DNA sequence.

#### Construction of the test plasmid:

[0339] Selection systems are the same as used in Example 1, *viz.* fusaric acid to select against cells expressing the tet gene and galactose killing by galT.K in a galE deleted host. First, in three genetic engineering steps, we replace: a) the ray gene in pEP1009 with the arc gene, and b) the target DNA sequences (both occurrences) with the arc operator. The resulting plasmid is our wild type control.

[0340] To replace ray with arc, the synthetic arc gene, shown in Table 201 and Table 202, is synthesized and ligated into pEP1009 that has been digested with BstEII and KpnI. Cells are transformed and colonies are screened for  $Tc^R$ . The plasmid is named pEP2000. Delta4 cells transformed with pEP2000 are  $Tc^R$  and  $Gal^S$  because pEP2000 lacks the ray gene.

[0341] To insert the arc operator into the neo promoter ( $P_{neo}$ ) for the tet gene in pEP2000, we digest pEP2000 with StuI and HindIII and ligate the purified backbone to annealed synthetic olig#430 and olig#432.



Arc operator and P<sub>neo</sub> that promotes tet

5                    5' | CCT|GCG|AAC|CGG|AAT|TGC|CAG|-  
 Olig #430 = 3' gga cgc ttg gcc tta acg gtc-  
                   | StuI |                    | -35 |

10 | CTG|GGG|CGC|CCT|CTG|GTA|AGG|TTG|-  
     gac ccc gcg gga gac cat tcc aac-  
                   | -10 |

15 | GGA|ATG|ATA|GAA|GCA|CTC|TAC|TAT|A                    3'-Olig#432  
     cct tac tat ctt cgt gag atg ata t tcg a 5'  
 20 | Arc operator | | Hind3 |

The plasmid is named pEP2001 and confers Fus<sup>R</sup>, Gal<sup>S</sup>, Ap<sup>R</sup> on delta4 cells.

25 [0342] To insert the arc operator into the amp promoter for the galT,K genes in pEP2001, we digest pEP2001 with Apal and XbaI and ligate the purified backbone to synthetic olig#416 and olig#417 that have been annealed in the standard way.

Arc operator and P<sub>amp</sub> that promotes galT,K

30                    5' | CTT|CTA|AAT|ACA|TTC|AAA|-  
                   Olig#417 3' c cgg gaa gat tta tgt aag ttt-  
                   | Apal |                    | -35 |

35 | TAT|GTA|TCC|GCT|CAT|GAG|ACA|ATA|ACC|-  
     ata cat agg cga gta ctc tgt tat tgg-  
                   | -10 |

40 | CTT|ATG|ATA|GAA|GCA|CTC|TAC|TAT| CGT                    3'Olig#416  
     gaa tac tat ctt cgt gag atg ata gca gat c 5'  
 45 | Arc Operator | | XbaI |

50 The plasmid is named pEP2002 and confers Gal<sup>R</sup>, Fus<sup>R</sup>, Ap<sup>R</sup> on delta4 cells. This plasmid is our wild type for work with polypeptides that are selected for binding to target DNA subsequences that are related to the arc operator.

Development of polypeptides that bind chimeric target DNA:

55

[0343] We now replace:

- a) the two occurrences of the arc operator with the first target sequence that is a hybrid of the arc operator and a

subsequence picked from HIV-1, and

b) the arc gene by a variegated pdbp gene.

5 [0344] A hybrid non-palindromic target sequence is used in this example because selection of a polypeptide using a palindromic or nearly palindromic target DNA subsequence is likely to isolate a novel dimeric DBP. The goal of this procedure is to isolate a polypeptide that binds DNA but that does not directly exploit the dyad symmetry of DNA. The binding is most likely in the major groove, but the present invention is not limited to polypeptides that bind in the major groove. The selections are performed using a non-symmetric target to avoid isolation of novel dimers that support two

10 symmetrically related copies of the original recognition elements.  
[0345] The non-variable regions of the HIV-1 genome, as listed in Example 1, were searched using a half operator from the arc operator as search sequence.

[0346] We sought subsequences in the non-variable sequences of the HIV-1 genome that match either half of the consensus P22 arc operator shown in Table 200. Subsequences that are closer to the start of transcription are preferred as targets because proteins binding to these subsequences will have greater effect on the transcription of the genes. No sequence was found that matched all six unambiguous bases; the subsequences at 1024, 1040, and 2387 (shown in Table 203) each have a single mismatch. Lower case letters in the "arcQ =" sequence indicate ambiguity in the P22 arc operator sequence. Lower case, bold, underscored letters in the HIV-1 subsequences indicate mismatch with the consensus arc operator. Two other subsequences, shown in Table 203, have one mismatch at one of the conserved bases and one mismatch with one of the ambiguous bases. The HIV-1 subsequence that starts at base 1024 is chosen as a target sequence. We replace the 3' ten bases of the arc operator with the 3' ten bases of this subsequence to produce the hybrid target sequence:

ATGATAGAAG[C]GCAACCCCTC . We insert this sequence into the promotor that regulates tet in pEP2002 by ligating dsDNA composed of an equimolar mixture of olig#440 and olig#442 into the StuI/HindIII site of pEP2002. Substitution of the arc operator by the arc-HIV-1 hybrid sequence relieves the repression by Arc. The construction is called pEP2003 and confers Tc<sup>R</sup>, Ap<sup>R</sup>, Gal<sup>S</sup> on delta4 cells.

#### First Target and P<sub>neo</sub> that promotes tet

30  
5' | CCT|GCG|AAC|CGG|AAT|TGC|CAG|-  
Olig#440 = 3' gga cgc ttg gcc tta acg gtc-  
| StuI | | -35 |

35  
| CTG|GGG|CGC|CCT|CTG|GTA|AGG|TTG|-  
gac ccc gcg gga gac cat tcc aac-  
40 | -10 |

45  
ATA ATA CAG TAg caa ccc tct = HIV 1024-1044  
| GGA|ATG|ATA|GAA|GCg|caa|ccc|tct|A 3'=Olig#442  
cct tac tat ctt cgC GTT GGG AGA t tcg a 5'  
| First Target | | Hind3 |

50  
[0347] The second instance of the target is engineered in like manner, using pEP2003 first digested with ApaI and XbaI and then ligated to annealed olig#444 and olig#446. The plasmid is called pEP2004 and confers Gal<sup>+</sup>, Tc<sup>R</sup>, Ap<sup>R</sup> on HB101 cells. The plasmid pEP2004 contains the first target sequence in both selectable genes and is ready for introduction of a variegated pdbp gene.

First Target and P<sub>amp</sub> that promotes galT.K

5'                    |CTT|CTA|AAT|ACA|TTC|AAA|  
 Olig#444    3'    c cgg gaa gat tta tgt|aag ttt|  
                   | AdAI |                    | -35 |

10

|TAT|GTA|TCC|GCT|CAT|GAG|ACA|ATA|ACC|CT-  
 ata cat agg cga gta ctc tgt tat tgg ga  
    | -10 |

20

T|ATG|ATA|GAA|GCg|caa|ccc|tct| CGT                    3'Olig#446  
 a tac tat ctt cgC GTT GGG AGA    gca gat c    5'  
 | First Target |                    | XbaI |

25

[0348]    The variegated DNA for a 36 amino acid polypeptide is shown in Table 204. This DNA encodes the first ten amino acids of P22 Arc followed by 26 amino acids chosen to be likely to form extended structures. In Table 204, we indicate variegation at one base by using a letter, other than A, C, G, or T, to represent a specific mixture of deoxynucleotide substrates. The range of amino acids encoded is written above the codon number:

35                    I|M  
    | 11 |  
    | ATs |

40 indicates that the first base is synthesized with A, the second base with T, and the third base with a mixture of C and G, and that the resulting DNA could encode amino acids I or M. That the parental protein has isoleucine at residue 11 is indicated by writing I first. Residues 22 and 23 are not variegated to provide a homologous overlap region so that olig#420 and olig#421 can be annealed. After olig#420 and olig#421 are annealed and extended with Klenow fragment and all four deoxynucleotide triphosphates, the DNA is digested with both BstEII and Bsu36I and ligated into pEP2004 that has also been digested with BstEII and Bsu36I. The ligated DNA, denoted vg1-pEP2004, is used to transform Delta4 cells. After an appropriate grow out in the presence of IPTG, the cells are selected with fusaric acid and galactose.

[0349]    By hypothesis, we recover ten colonies that are Gal<sup>R</sup> and Fus<sup>R</sup>. We sequence the plasmid DNA from each of these colonies. A hypothetical DBP amino acid sequence from one of these colonies is shown in Table 205.

50 [0350]    Comparison of the amino-acid sequences of different isolates may provide useful information on which residues play crucial roles in DNA binding. Should a residue contain the same amino acid in most or all isolates, we might infer that the selected amino acids is preferred for binding to the target sequence. Because we do not know that all of the isolates bind in the same manner, this inference must be considered as tentative. Residues closer to the unvaried section that have repetitive isolates containing the same amino acid are more informative than residues farther away.

55 [0351]    In a second round of Diffuse Mutagenesis, we vary the codons shown in Table 206. Residues 1 through 10 are not varied because these provide the best match for the first ten bases of the target. Residues 19, 20, and 21 are not varied so that the synthetic oligonucleotides can be annealed. The two-way variations at residues 11 through 18 and 23 through 36 all allow the selected amino acid to be present, but also allow an as-yet-untested amino acid to appear.

It is desirable to introduce as much variegation as the genetic engineering and selection methods can tolerate without risk that the parental DBP sequence will fall below detectable level. Having picked three residues for the homologous overlap, we have only 23 amino acids to vary. Thus residue 22 is varied through four possibilities instead of only two. Residue 22 was chosen for four-way variegation because it is next to the unvaried residues. We use pEP2004 as the backbone, and ligate DNA prepared with Klenow fragment from oligonucleotides #423 and #424 (Table 206) to the BstII and Bsu36I sites. The resulting population of plasmids containing the variegated DNA is denoted vg2-pEP2004.

[0352] Table 207 shows the amino acid sequence obtained from a hypothetical isolate bearing a DBP gene specifying a polypeptide with improved affinity for the target. Changes in amino acid sequence are observed at ten positions. Comparisons of the sequences from several such isolates as well as those obtained in the first round of mutagenesis can be used to locate residues providing significant DNA-binding energy.

[0353] Having established some affinity for the target, we now seek to optimize binding via a more focused mutagenesis procedure. Table 208 shows a third variegation in which twelve residues in the variable region are varied through four amino acids in such a way that the previously selected amino acids may occur. Again, pEP2004 is used as backbone and synthetic DNA having cohesive ends is prepared from olig#325 and olig#327. The plasmid is denoted vg3-pEP2004. In subsequent variegation, we would vary other residues through four amino acids at one time. By hypothesis, we select the polypeptide shown in Table 209 that has high specific affinity for the first target; now we can:

a) replace both occurrences of the first target by a second target, i.e. the intact HIV-1 subsequence (1024-1044), and

b) use the selected polypeptide as the parental DBP to generate a variegated population of polypeptides from which we select one or more that bind to the second target.

Because the second target differs from the first in the region thought to be bound by residues 1 through 10 of the parental DBP, we concentrate our variegation within these residues for the first several rounds of variegation and selection.

[0354] We replace the target DNA sequence in the neo promoter for tet in pEP2002 with ds DNA comprising synthetic olig#450 and olig#452 at the StuI/HindIII site. The new plasmid is named pEP2010 and confers Tc<sup>R</sup> on delta4 cells.

### Second Target and P<sub>neo</sub> that promotes tet

```

5'   | CCT|GCG|AAC|CGG|AAT|TGC|CAG|-
Olig#450 = 3'   gga cgc ttg gcc tta acg gtc-
           | StuI |           | -35 |

```

```

| CTG|GGG|CGC|CCT|CTG|GTA|AGG|TTG|GG-
gac ccc gcg gga gac cat tcc aac cc-
           | -10 |

```

```

ATA ATA CAG TAg caa ccc tct = HIV 1024-1044
A|ATa|ATA|cAg|taG|caa|ccc|tct|A      3'Olig#452
t taT tat GtC ATC GTT GGG AGA t tcg a 5'
| Second Target | | Hind3 |

```

[0355] We replace the target in the amp promoter for galT<sub>K</sub> of pEP2010 with synthetic olig#454 and olig#456 between ApaI and XbaI sites. The new plasmid is named pEP2011 and confers Gal<sup>+</sup> on HB101. pEP2011 contains the second target in both selectable genes and is ready for introduction of a variegated gdbp gene and selection of cells expressing polypeptides that can selectively bind the target DNA subsequence.

Second Target and P<sub>amp</sub> that promotes galT.K

5

Olig#454

5' |CTT|CTA|AAT|ACA|TTC|AAA|

3' c cgg gaa gat tta tgt|aag ttt|

|ApaI| | -35 |

15 |TAT|GTA|TCC|GCT|CAT|GAG|ACA|ATA|ACC|CT  
ata cat agg cga gta ctc tgt tat tgg ga  
| -10 - |

20 ATA ATA CAG Tag caa ccc tct = HIV 1024-1044  
T|ATa|ATA|cAg|tag|caa|ccc|tct| CGT 3'Olig#456  
25 a taT tat GtC ATc GTT GGG AGA gca gat c 5'  
| Second Target | | XbaI |

30 **[0356]** Variation of the first eleven residues of the potential DNA-binding polypeptide is illustrated in Table 210. Double-stranded DNA having appropriate cohesive ends is prepared from olig#460 and olig#461, Klenow fragment, BstEII, and Bsu36I. This DNA is ligated into similarly digested backbone DNA from pEP2011; the resulting plasmid is denoted vg1-pEP2011. Delta4 cells are transformed and selected with fusaric acid and galactose. Table 211 shows the sequence of a 37 amino-acid polypeptide isolated from cells exhibiting the DBP<sup>+</sup> phenotypes by the above hypothetical  
35 selection. The sequence shown in Table 211 is hypothetical and is given by way of example. This example must not be construed as a prediction that this sequence has specific affinity for the target or any other DNA sequence. Further var-  
iegation (vg2, vg3,...) of this peptide and selection for binding to Target#2 will be needed to obtain a peptide of high spe-  
cificity and affinity for Target#2.

[0357] We anticipate that Successful DBP production will take more than three or four cycles of variegation and selection, perhaps 10 or 15. We anticipate that initial phases will require careful adjustment of the selective agents and IPTG because the level of repression afforded by the best polypeptide may be quite low. As stated, we expect that biophysical methods, such as X-ray diffraction or NMR, applied to complexes of DNA and polypeptide will yield important indications of how to hasten the forced evolution.

**[0358]** The length of the polypeptide in the example may not be optimal; longer or shorter polypeptides may be needed. It may be necessary to bias the amino acid composition more toward basic amino acids in initial phases to obtain some non-specific DNA binding. Inclusion of numerous aromatic amino acids (W,F,Y,H) may be helpful or necessary.

Other strategies to obtain polypeptides that bind sequence-specifically are illustrated in examples 3, 4, and 5.

50 **Example 3**

**[0360]** We present a second example of the application of our selection method applied to the generation of asymmetric DBPs. A possible problem with making and using DNA-binding polypeptides, is that the polypeptides may be degraded in the cell before they can bind to DNA. That polypeptides can bind to DNA is evident from the information on sequence-specific binding of oligopeptides such as Hoechst 33258. Polypeptides composed of the 20 common natural amino acids contain all the needed groups to bind DNA sequence-specifically. These are obtained by an efficient method to sort out the sequences that bind to the chosen target from the ones that do not. To overcome the tendency

of the cells to degrade polypeptides, we will attach a domain of protein to the variegated polypeptide as a custodian. The first example of a custodial domain presented is residues 20-83 of barley chymotrypsin inhibitor.

[0361] The strategy is to fuse a polypeptide sequence to a stable protein, assuming that the polypeptide will fold up on the stable domain and be relatively more protected from proteases than the free polypeptide would be. If the domain is stable enough, then the polypeptide tail will form a make-shift structure on the surface of the stable domain, but when the DNA is present, the polypeptide tail will quickly (a few milliseconds) abandon its former protector and bind the DNA. The barley chymotrypsin inhibitor (BCI-2) is chosen because it is a very stable domain that does not depend on disulfide bonds for stability. We could attach the variegated tail at either end of BCI-2. A preferred order of amino acid residues in the chimeric polypeptide is: a) methionine to initiate translation, b) BCI-2 residues 20-83, c) a two residue linker, d) the first ten residues of Arc, and e) twenty-four residues that are varied over two amino acids at each residue. The linker consists of G-K. Glycine is chosen to impart flexibility. Lysine is included to provide the potentially important free amino group formerly available at the amino terminus of the Arc protein. The first target is the same as the first target of Example 2.

[0362] Table 300 shows the sequence of a gene encoding the required sequence. The ambiguity of the genetic code has been resolved to create restriction sites for enzymes that do not cut pEP1009 outside the arc gene. This gene could be synthesized in several ways, including the method illustrated in Table 301 involving ligation of oligonucleotides 470-479. Plasmid pEP3000 is derived from pEP2004 by replacement of the arc gene with the sequence shown in Table 300 by any appropriate method.

[0363] Table 302 illustrates variegated olig#480 and olig#481 that are annealed and introduced into the C12-arc(1-10) gene between PpuMI and KpnI to produce the plasmid population vg1-pEP3000. Cells transformed with vg1-pEP3000 are selected with fusaric acid and galactose in the presence of IPTG. Further variegation (vg2, vg3, ...) will be required to obtain a polypeptide sequence having acceptably high specificity and affinity for Target#1.

#### Example 4

[0364] We present a second strategy involving a polypeptide chain attached to a custodial domain. In this strategy, the custodial domain contains a DNA-recognizing element that will be exploited to obtain quicker convergence of the forced evolution.

[0365] The three alpha helices of Cro fold on each other. It has not been observed that these helices fold by themselves, but no efforts in this direction have been reported. We will attach a variegated segment of 24 residues to residue 35 of Cro (H35 is the last residue of alpha 3). The target will be picked to contain a good approximation to the half  $O_{R3}$  site at one end but no constraint is placed on the bases corresponding to the dyad-related other half of  $O_{R3}$ . A sequence that departs widely from the  $O_{R3}$  sequence is actually preferred, because this discourages selection of a novel dimeric molecule. We assume that alpha-3 forms and binds to the same four or five bases that it binds in  $O_{R3}$  and that a polypeptide segment attached to the carboxy terminus of alpha-3 can continue along the major groove. We attach 24 amino acids of polypeptide immediately after the last residue of alpha-3, wherein the polypeptide is chosen: a) to have more positive charge than negative charge, b) to have beta chain predominate, c) to have some aromatic groups, and d) to have some H-bonding groups, produces a population that is then cloned and host cells are selected for expression of a polypeptide that binds preferentially to the target sequence.

[0366] We first construct a hybrid target sequence (Target #3) containing one  $O_{R3}$  half-site fused to a portion of the final target. This hybrid target DNA subsequence is inserted into the selectable genes in the same manner as the arc operator was inserted in Example 2. We then follow the same procedure to vary the 24 residues; first we vary twenty-four residues, using two possible amino acids at each residue. We carry out two or more cycles of such diffuse variegation. Then we vary 12 residues, using 4 possible amino acids at each residue. We do two or more iterations of this process so that all residues are varied at least once.

[0367] We have now generated one or more DBPs that bind well to one half of the final target sequence. Next we generate binding to the other half of the final target. First we replace both instances of Target #3 with the final target sequence, target #4. We then vary the alpha helix 3 and the surface of the hypothesized domain formed by helices 1-3 to optimize binding to final target sequence.

[0368] A search of the non-variable regions of the HIV-1 genome reveals that bases 624-640 (aATCCTAGCAGTGGCG) contain a good match to one half of  $O_{R3}$ , as shown in Table 400. As first target of this example, we choose TATCCCTAGCAGTGGCG, denoted Target#3, that has one half of  $O_{R3}$  and nine bases from HIV-1. Once a sequence is obtained that binds Target#3, we replace Target#3 by Target#4 = HIV 624-640 and variegate the recognition helices taken from Cro.

[0369] To engineer Target#3 into  $P_{neo}$  that regulates tet, plasmid pEP2002 is digested with StuI and HindIII and the purified backbone is ligated to an annealed, equimolar mixture of olig#490 and olig#492. Delta4 cells are transformed and selected with Tc; replacement of the arc operator relieves the repression by Arc. Plasmid DNA from Tc<sup>R</sup> colonies is sequenced to confirm the construction; the construction is called pEP4000.

Target #3 and P<sub>neo</sub> that promotes tet

5' | CCT|GCG|AAC|CGG|AAT|TGC|CAG|-  
 Olig#490 = 3' gga cgc ttg gcc tta acg gtc-  
                   | StuI |                   | -35 |

10

15 | CTG|GGG|CGC|CCT|CTG|GTA|AGG|TTG|GG-  
     gac ccc gcg gga gac cat tcc aac cc-  
                                   | -10 |

20       aAT CtC TAG CAG TGG CG = HIV 624-640  
 A|TAT|CCC|TAG|CAG|TGG|CGA       3'Olig#492  
 t ata ggg atc gtc acc gct tcg a 5'

25 | \_\_\_\_\_ Target #3 \_\_\_\_\_ | |Hind3| |

[0370] We engineer the second instance of the target, in like manner, into P<sub>amp</sub> for galT.K, using ApaI and XbaI to digest pEP4000 and olig#494 and olig#496. HB101 cells (galK) are transformed and are selected for ability to grow on galactose as sole carbon source. Plasmid DNA from Gal<sup>+</sup> colonies is sequenced in the region of the insert to confirm the construction. The plasmid is called pEP4001.

Target #3 and P<sub>amp</sub> that promotes galT.K

35                   5'           | CTT|CTA|AAT|ACA|TTC|AAA|  
 Olig#494       3' c cgg gaa gat tta tgt|aag ttt|  
                                   | ApaI |                   | -35 |

40

45 | TAT|GTA|TCC|GCT|CAT|GAG|ACA|ATA|ACC|  
     ata cat agg cga gta ctc tgt tat tgg  
                                   | -10 |

50       | CTT|TAT|CCC|TAG|CAG|TGG|CG CGT       3'Olig#496  
     gaa ata ggg atc gtc acc gc gca gat c 5'

55 | \_\_\_\_\_ Target #3 \_\_\_\_\_ | |XbaI| |

[0371] A gene fragment encoding the first two helices of Cro is shown in Table 401. Olig#483 and olig#484 are syn-

thesized and extended in the standard manner and the DNA is digested with BstEII and KpnI. This DNA is ligated to backbone from pEP4001 that has been digested with BstEII and KpnI; the resulting plasmid, denoted pEP4002, contains the Target#3 subsequence in both selectable genes and is ready for introduction of a variegated pdbp gene between BglII and KpnI. Table 402 shows a piece of vgDNA prepared to be inserted into the BglII-KpnI sites of pEP4002. Table 403 shows the result of a selection of delta4 cells, transformed with vg1-pEP4002, with fusaric acid and galactose in the presence of IPTG. Additional cycles of variegation of residues 36-61 are carried out in such a way that the amino acid selected at the previous cycle is included. After several cycles in which 22-24 residues are varied through two possible amino acids, we choose 10-13 amino acids and vary them through four possibilities.

[0372] Once reasonably strong binding to Target#3 is obtained, we replace Target#3 with Target#4 and vary the residues in helix 3 (residues 26-35) and, to a lesser extent, helix 2 (residues 16-23).

#### Example 5

[0373] We disclose here a method of engineering a polypeptide extension onto the amino terminus of P22 Arc, a natural DBP, so that the novel DBP develops asymmetric DNA-binding specificity for a subsequence found in the HIV-1 genome. Others have observed that loss of arms from natural DBPs may cause loss of binding specificity and affinity (PABO82a and ELIA85), but none, to our knowledge, have suggested adding arms to natural DBPs in order to enhance or alter specificity or affinity. The new construction is denoted a "polypeptide extension DBP"; the gene is denoted ped and the proteins are denoted Ped. Wild-type Arc forms dimers and binds to a partially palindromic operator. We will generate a sequence of DBPs descendent from Arc. Early members of this family will form dimers, but will have sufficient binding area such that asymmetric targets will be bound. In final stages of the development, proteins that do not dimerize will be engineered.

[0374] Table 200 shows the symmetric consensus of left and right halves of the P22 arcO operator, arcO. Table 500a shows a schematic representation of the model for binding of Arc to arcO that is supported by genetic and biochemical data (VERS87b). Arc is thought to bind B-DNA in such a way that residues 1-10 are extended and the amino terminus of each monomer contacts the outer bases of the 21 bp operator (RT Sauer, public talk at MIT, 15 September 1987).

[0375] Arc is preferred because: a) one end of the polypeptide chain is thought to contact the DNA at the exterior edge of the operator, and b) Arc is quite small so that genetic engineering is facilitated. P22 Mnt is also a good candidate for this strategy because it is thought that the amino terminal six residues contact the mntO operator, mntO, in substantially the same manner as Arc contacts arcO. Mnt has significant (40%) sequence similarity to Arc (VERS87a). Mnt forms tetramers in solution and it is thought that the tetramers bind DNA while other forms do not. When the mnt gene is progressively deleted from the 3' end to encode truncated proteins, it is observed that proteins lacking K79 and subsequent residues have lowered affinity for mntO and that proteins lacking Y78 and subsequent residues can not form tetramers and do not bind DNA sequence-specifically (KNIG88). Some truncated Mnt proteins of 77 or fewer residues form dimers, but these dimers do not present the DNA-recognizing elements in such a way that DNA can be bound. Arc is preferred over Mnt because Arc is smaller and because Arc acts as a dimer.

[0376] Other natural DBPs that have DNA-recognizing segments thought to interact with DNA in an extended conformation (referred to as arms or tails) and thought to contact the central part of the operator, such as  $\lambda$  Cro or  $\lambda$  cl repressor, are less useful. For these proteins to be lengthened enough to contact DNA outside the original operator, several residues would be needed to span the space between the central bases contacted by the existing terminal residues and the exterior edge of the operator.

[0377] Table 500a illustrates interaction of Arc dimers with arcO; the two "C"s of Arc represent the place, near residue F10, at which the polypeptide chain ceases to make direct contact with the DNA and folds back on itself to form a globular domain, as shown in Table 500b and Table 500c. Which of these alternative possibilities actually occurs has not been reported. Our strategy is compatible, with some alterations, with either structure. In Table 500b, each set of residues 1-10 makes contact with a domain composed of residues 11-57 of the same polypeptide chain; the dimer contacts are near the carboxy terminus. Table 500c shows an alternative interaction in which residues 1-10 of one polypeptide chains interact with residues 11-57 of the other polypeptide chain; the dimer contacts occur shortly after residue 10. The similarity of sequences of Arc and Mnt, the demonstration of function of DNA-recognizing segments transferred from Arc to Mnt (RT Sauer, public talk at MIT, 15 September 1987 and Knight and Sauer cited in VERS86b), and the behavior of Mnt on truncation suggest that Table 500b is the correct general structure for Arc, but the structure diagrammed in Table 500c is also possible.

[0378] Table 501 shows the four sites at which one of the consensus arc half operators comes within one base of matching ten bases (six unambiguous and four having two-fold ambiguity) in the non-variable segments of HIV-1 DNA sequence, as listed in Example 1. The symbol "@" marks base pairs that vary among different strains of HIV-1. Because we intend to extend Arc from its amino terminus, we seek subsequences of HIV-1 that: a) match one of the arc half operators, and b) have non-variable sequences located so that an amino-terminal extension of the Arc protein will interact with non-variable DNA. The subsequences 1024-1033 and 4676-4685 meet this requirement while the sub-



sequences at 1040-1049 and 2387-2396 do not. In the case of 1040-1049, the amino-terminal extension would proceed in the 3' direction of the strand shown and would reach variable DNA after two base pairs. For 2387-2396, variable sequence is reached at once. The subsequence 1024-1033 is preferred over the subsequence 4676-4685 because it is much closer to the beginning of transcription of HIV so that binding of a protein at this site will have a much greater effect on transcription. In the remainder of this example, positions within the target DNA sequence will be given the number of the corresponding base in HIV-1. Base A<sub>1034</sub> of HIV-1 is aligned with the central base of arcQ.

[0379] HIV 1024-1044 has only three bases in each half that are palindromically related to bases in the other half by rotation about base pair 1034: A<sub>1024</sub>/T<sub>1044</sub>, A<sub>1026</sub>/T<sub>1042</sub>, and G<sub>1032</sub>/C<sub>1036</sub>. The latter two base pairs correspond to positions in arcQ that are not palindromically related. Five of the six palindromically related bases of arcQ correspond to non-palindromically related bases in HIV 1024-1044. Thus no dimeric protein derived from Arc is likely to bind HIV 1016-1046 if symmetric changes are made only in the residues 1-10 (or in any other set of residues originally found in Arc). Our strategy is to add, in stages, eleven variegated residues at the amino terminus and to select for specific binding to a progression of targets, the final target of the progression being bases 1016-1037 of HIV-1. Because the region of protein-DNA interaction is increased beyond that inferred for wild-type Arc-arcQ complexes, unfavorable contacts in bases aligned with the right half of arcQ can be compensated by favorable contacts of the polypeptide extension with bases 1016-1023. The penultimate selection isolates a dimeric protein that binds to the HIV-1 target 1016-1037; the ultimate selection isolates a protein that does not dimerize and binds to the same target.

[0380] Table 502 shows a progression of target sequences that leads from wild-type arcQ to HIV 1016-1037. It is emphasized that finding a subsequence of HIV-1 that has high similarity to one half of arcQ is not necessary; rather, use of this similarity reduces the number of steps needed to change a sequence that is highly similar to arcQ into one that is highly similar or identical to an HIV-1 subsequence. Reducing the number of steps is useful, because, for each change in target, we must: a) construct plasmids bearing selectable genes that include the target sequence in the promoter region, b) construct a variegated population of ped genes, and c) select cells transformed with plasmids carrying the variegated population of ped genes for DBP<sup>+</sup> phenotype.

[0381] In sections (a), (c), (e), and (g) of Table 502, bases in the targets are in upper case if they match HIV 1016-1046 and are underscored if they match the wild-type arcQ sequence.

[0382] We construct a series of plasmids, each plasmid containing one of the target sequences in the promoter region of each of the selectable genes. For each target, we variegate the ped gene and select cells for phenotypes dependent on functional DBPs. For each target, several rounds of variegation and selection may be required. We anticipate that a plurality of proteins will be obtained from independent isolates by selection for binding to one target. We pick the protein that shows the strongest *in vitro* binding to short DNA segments containing the target as the parental Ped to the next round of variegation and selection. Genetic methods, such as generation of point mutations in the ped gene or in the target and selection for function or non-function of Ped can be used to determine associations between particular bases and particular residues (VERS86b).

[0383] Once a Ped with specific binding for the target is obtained, it may be useful to determine a 3D structure of the Ped-DNA complex by X-ray diffraction or other suitable means. Such a structure would provide great help in choosing residues to vary to improve binding to a given target or to an altered target.

[0384] We initiate development of a polypeptide extension DBP having affinity for HIV 1016-1037 by generating a variegated population of Peds and selecting for binding to the first target. Table 502a shows the first target which we designed to have identity to arcQ in the left half, but to have a mismatch (arcQ vs. target) at A<sub>1038</sub> (which is C in the corresponding position in the right half of arcQ and is palindromically related to a G in the left half); the rationale is as follows. Vershon *et al.* (VERS87b) report that chemical modification with dimethyl sulfate of the wild-type CG at this location interferes mildly with binding of Arc and that this location is strongly protected from modification by dimethylsulfate if Arc is bound to the operator. Thus we expect a mismatch between wild-type arcQ and the first target at A<sub>1038</sub> to make wild-type Arc bind poorly. Binding can be restored, however, by favorable contacts to bases 1021-1023 by the amino-terminal extension.

[0385] An alternative first target would have C<sub>1038</sub>, as does arcQ at the corresponding location, and A<sub>1041</sub>, unlike arcQ or HIV-1. Vershon *et al.* (VERS87b) report that methylation of the corresponding CG base pair strongly interferes with binding of Arc. Thus, changing the base that corresponds to HIV 1041 should have a strong effect on binding of Arc to the alternative target.

[0386] In the first variegation step, we extend Arc by five variegated residues at the amino terminal. Since five residues can contact no more than three bases in a sequence-specific manner, we limit the extent of the target to those bases that correspond to HIV 1021-1044. Inclusion of bases corresponding to HIV 1016-1020 at this initial stage might position the target too far downstream from the promoters of the selectable genes to allow strong repression of these promoters. Once a Ped displaying binding to bases corresponding to 1021-1044 has been isolated, we can introduce a greater length of the HIV-1 sequence into the left side of the target without concern that the Ped will bind too far downstream from the promoter of the selectable genes to block transcription. Furthermore, once binding by the amino terminal extension has been established, we can, in a stepwise manner, remove the right half of arcQ from the target.

thereby forcing more asymmetric binding to the left half of arcQ and the bases upstream of 1024.

[0387] The first target is engineered into both selectable genes as in Example 2. We use olig#501 and olig#502, shown in Table 503, to introduce the first target downstream of  $P_{neo}$  that promotes tet, replacing arcQ in pEP2002; the resulting plasmid is called pEP5000. From pEP5000, we use olig#503 and olig#504 to construct pEP5010 in which the first target replaces arcQ downstream of  $P_{amp}$  that promotes galT<sub>L</sub>K.

[0388] Table 502b shows schematically how the amino terminal residues align to the first target; the five residue extension is unlikely to contact more than 3 base pairs upstream from base 1024. The alteration in the right half operator prevents tight binding unless the additional residues make favorable interactions upstream of 1024. Care is taken in designing the two instances of the target that the downstream boundaries are different, AAG in  $P_{neo}$  and CGT in  $P_{amp}$ . Thus, for the novel DBP to bind specifically to both instances of the target, it must recognize the common sequence upstream of base 1024.

[0389] An initial variegated ped is constructed using olig#605, as shown in Table 504, and comprises: a) a methionine codon to initiate translation, b) five variegated codons that each allow all twenty possible amino acids, and c) the Arc sequence from 101 to 157. (Because we are constructing a polypeptide extension at the amino terminus, we have added 100 to the residue numbers within Arc so that Arc residue 1 is designated 101.) This variegated segment of DNA comprises  $(2^5)^5 = 2^{25} = 3.2 \times 10^7$  different DNA sequences and encodes  $20^5 = 3.2 \times 10^6$  different protein sequences; with the given technical capabilities, we can detect each of the possible protein sequences. The 3' terminal 20 bases of olig#605 are palindromically related so that each synthetic oligonucleotide primes itself for extension with Klenow enzyme. The DNA is then digested with Bsu36I and BstEII and is ligated to the backbone of appropriately digested pEP5010 which bears the first target in each selectable gene. Transformed delta4 cells are selected for  $Fus^R$   $Gal^R$  at low, medium, and high concentrations of IPTG, the inducer of the lacUV5 promoter that regulates ped. Because the first target is quite similar to arcQ, we anticipate that a functional Ped will be isolated with low-level induction of the ped gene with IPTG.

[0390] More than one round of variegation and selection may be required to obtain a Ped with sufficient affinity and specificity for the first target. Function of a Ped is judged in comparison to the protection afforded by wild-type Arc in cells bearing pEP2002. Specifically, strength of Ped binding is measured by the IPTG concentration at which 50% of cells survive selection with a constant concentration of galactose or fusaric acid, chosen as a standard for this purpose. A Ped is deemed acceptable if it can protect cells against the standard concentrations of galactose and fusaric acid, administered in separate tests, with an IPTG concentration of  $5 \times 10^{-4}$  M. Preferably, a Ped can protect cells against the standard concentrations of galactose and fusaric acid, tested separately, with no more than ten times the concentration of IPTG needed by pEP2002-bearing cells. Variegation of residues 101, 102, and others may be needed. We anticipate that a plurality of independent functional Peds will be isolated; we discriminate among these by measuring *in vitro* binding to DNA oligonucleotides that contain the target sequence. The amino-acid sequences of different isolates are compared; residues that always contain only one or a few kinds of amino acids are likely to be involved in sequence-specific DNA binding. Table 505 shows a hypothetical isolate, Ped-6, that binds the first target.

[0391] Table 502c shows the changes between the first target and the second target. Three changes are made left of center to make the target more like HIV 1016-1042. Only the change  $G_{1030} \rightarrow C$  affects a base that is palindromically related in arcQ. One change is made right of center that makes the target more like HIV 1016-1042, less like arcQ, and less palindromically symmetric. Furthermore, the target is shortened on the right by two bases so that selection isolates proteins that bind asymmetrically to the left side of the target. Starting with pEP2002, we introduce, in two genetic engineering steps that use olig#541, olig#542, olig#543 and olig#544 (Table 506), the second target (in place of arcQ) into the promoter region of each selectable gene; the resulting plasmid is denoted pEP5020.

[0392] Table 507 shows a variegated sequence that is ligated into pEP5020 between BstEII and Bsu36I. Variegated codons are shown in the same way as in Table 204.

[0393] Table 502d illustrates that residues 100-110 of Ped-6 contact the bases of the second target that differ from the first target. Accordingly, residues 1 and 96-99 of Ped are not variegated in the DNA shown in Table 507; rather, residues 100-110 are each varied through four possibilities, always including the amino acid previously present at that residue. This generates  $4^{11} = 2^{22} = \text{approx. } 4 \times 10^6$  different DNA and protein sequences. Selection of transformed delta4 cells for  $Fus^R$   $Gal^R$  and screening by *in vitro* DNA binding yields, by hypothesis, a plasmid coding on expression for the protein Ped-6-2, illustrated in Table 508.

[0394] An alternative to the variegation shown in Table 507 is one in which we vary residues 101-105, 108, and 110 through eight possibilities each, yielding  $2.0 \times 10^6$  DNA and protein sequences. These residues, except M101, are indicated to be in contact with the operator. M101 has been altered by the attachment of the polypeptide extension and thus should be altered. After variegation of the listed residues and selection, further variegation should include some variegation of residues 96-103 because changes in the listed residues may change the context within which residues 96-103 contact the DNA.

[0395] More than one round of variegation and selection may be required to obtain a Ped having sufficient affinity and specificity for the second target.

[0396] Table 502e shows the changes from the second target to the third, which comprise: a) inclusion of bases 1018-1020, b) one change to the left of the 21 bp *arcQ* region, c) two changes at the center of the *arcQ* region, d) two changes left of center, and e) removal of bases 1041 and 1042. All of these changes make the third target less symmetric and more like HIV 1016-1040. The third target is introduced into each of the selectable genes in the same manner as the second target. The resulting plasmid, obtained in two genetic engineering steps, is denoted pEP5030. Table 502f shows that residues 96-110 are all potential sites to alter the specificity and affinity of DBPs derived from Ped-6-2. Thus, in Table 510, we illustrate a segment of variegated DNA that comprises  $2^{20} = 10^6$  DNA sequences and encoding on expression  $10^6$  protein sequences having ten residues varied through two possibilities and five residues through four possibilities. The DNA is then digested with *Bst*EI and *Bsu*36I and ligated into pEP5030. Transformed *delta4* cells are selected for *Fus*<sup>R</sup> *Gal*<sup>R</sup>. By hypothesis, we isolate a plasmid, denoted pEP5031, that codes on expression for the protein Ped-6-2-5 shown in Table 509.

[0397] Table 502g shows the changes between the third and fourth targets. The changes are: a) inclusion of bases 1016-1017, b) two changes right of center, and c) removal of bases 1038-1040. The initial variegation to be selected using the fourth target consists of an extension of six residues at the amino terminus of Ped-6-2-5, shown in Table 511. In iterative steps of forced evolution of proteins, one should not produce a number of different DNA sequences greater than the number of independent transformants that one can obtain (about  $10^8$  with current technology). Because there are no residues corresponding to 90-95 in the parental DBP (Ped-6-2-5), the first variegation and selection with the fourth target is a non-iterative step and it is permissible to produce  $10^{10}$  DNA sequences and  $6.4 \times 10^7$  protein sequences. In subsequent iterative rounds of variegation, the number of variants is, preferably, limited to a fraction, e.g., 10%, of the number of independent transformants that can be generated and subjected to selection. A protein, illustrated in Table 512 and denoted Ped-6-2-5-2, is isolated, by hypothesis, through selection of a variegated population of transformed cells for *Fus*<sup>R</sup> *Gal*<sup>R</sup>.

[0398] Ped-6-2-5-2 binds specifically to HIV 1016-1037 as a dimer. HIV 1016-1037 has no palindromic symmetry. Binding to an asymmetric DNA sequence by a dimeric protein is possible because the Ped-6-2-5-2 dimer has more recognition elements than wild-type P22 Arc dimer and so can bind even though nearly half of the right half of *arcQ* has been removed from the target. Ped-6-2-5-2 is useful as is; nevertheless, obtaining a monomeric protein may have advantages, including: a) higher affinity for the target because suboptimal interactions are eliminated, and b) lower molecular weight. Obtaining a functional monomeric Ped is easiest if Arc dimers interact in the manner shown in Table 500b. We use the following steps to isolate a protein that binds specifically to HIV 1016-1037 as a monomer.

[0399] Ped-6-2-5-2 is the parental DBP from which we derive the monomeric DBP. The route taken from a palindromically symmetric *arcQ* sequence to an asymmetric HIV sequence was designed to select for binding to the left half of the original *arc* operator.

[0400] Proteins that do not dimerize, but that bind specifically to the fourth target can be generated in several ways. Because the 3D structure of Arc is still unknown, we can not use Structure-Directed Mutagenesis to pick residues to vary to eliminate dimerization. One way to obtain monomeric proteins is to use diffuse mutagenesis to vary all residues from 111 to 157 and select for proteins that can bind the target sequence. Another strategy is to synthesize the *ped* gene in such a way that numerous stop codons are introduced. This causes a population of progressively truncated proteins to be expressed. Table 513 shows a segment of variegated DNA that spans the *Bgl*II to *Kpn*I sites of the *arc* gene used throughout this example. This segment is synthesized with suitable spacer sequences on the 5' end. The extra "T" at the 3' end allows two such chains to prime each other for extension with Klenow enzyme. The ratios of bases in the variegated positions are picked so that each varied codon encodes about 35% of polypeptides to terminate at that position. Since we intend to determine how much the protein can be shortened and remain functional, we begin by replacing codon 153 with stop. Since 15 residues are varied, only about 0.3 % of chains will continue to stop codon 153 without one or more stop codons. All the intermediate length chains will be present in the selection in detectable amount. *delta4* cells transformed with pEP5030 containing this vgDNA are selected for *Fus*<sup>R</sup> *Gal*<sup>R</sup>. Because each variegated codon causes translation termination in about 35% of the genes in the variegated population, shorter coding regions are more abundant than longer ones. Thus, the shortest gene that encodes a functional repressor will be the most abundant gene selected. Plasmid DNA from a number of independent selected colonies is sequenced. The dimerization properties of several functional DBPs are tested *in vitro* and the sequence of the shortest monomeric protein is retained for use and further study.

[0401] In this manner, we generate a protein that binds monomerically to a DNA sequence that has no palindromic symmetry.

#### Example 6

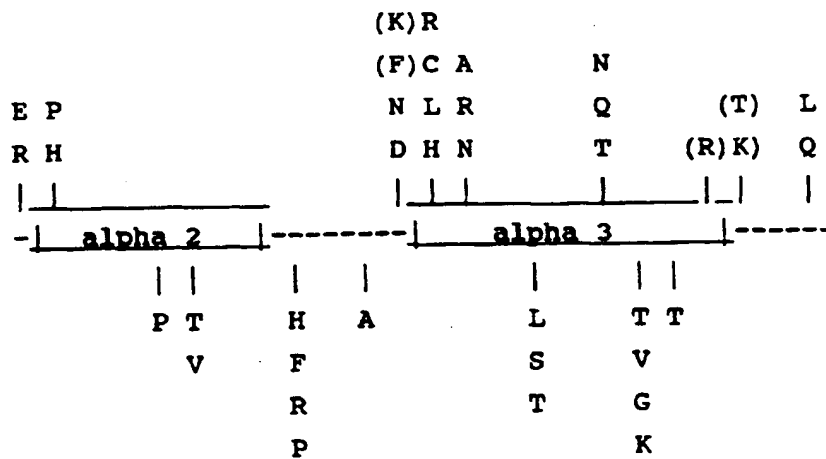
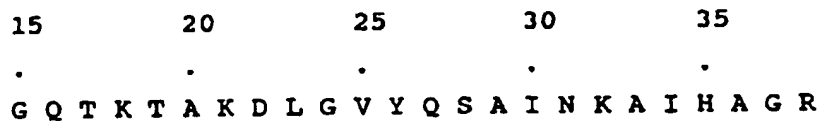
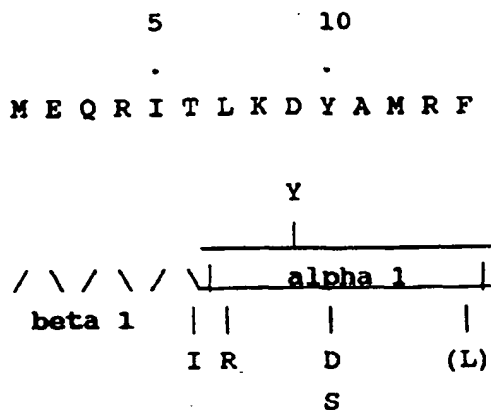
[0402] We illustrate here the fusion of two known DNA-binding domains to form a novel DNA-binding protein that recognizes an asymmetric target sequence. The progression of targets is the same as shown in Table 502 (Example 5). The amino-acid sequence of the initial DBP is illustrated in Table 600 and comprises the third zinc-finger domain

from the product of the *Drosophila* *kr* gene (ROSE86), a short linker, and P22 Arc. The linker consists of three residues that are picked to allow: a) some flexibility between the two domains, and b) introduction of a *Kpn*I site. The polypeptide linker should not allow excessive flexibility because this would reduce the specificity of the DBP.

[0403] The primary set of residues to vary to alter the DNA-binding are marked with asterisks. Those in the zinc finger were picked by reference to the model of Gibson *et al.* (GIBS88); all residues having outward-directed side groups (except those directed upward from the beta strands) were picked. Residues 101-110 (1-10 of Arc) were also picked to be in the primary set. Other residues within the Arc sequence may be varied. For each target in the progression, we initially choose for variegation residues in the primary set that are most likely to abut that part of the target most recently changed. For example, for the first target, we begin by varying residues 21, 24, 25, 28, and 29, each through all twenty amino acids. After one or more rounds of variegation and selection, other residues in the primary and secondary set are varied.

[0404] Other zinc-finger domains, such as those tabulated by Gibson *et al.* (GIBS88), are potential binding domains. Other proteins with known DNA binding, such as 434 Cro, may be used in place of Arc. Multiple zinc fingers could be added, stepwise, to obtain higher levels of specificity and affinity.

## TABLES

Table 1 MISSENSE MUTATIONS IN  $\lambda$  CRO

10

15.

20

25

30      **Notes:**

35

40

45

50

55

EP 0 452 413 B1

Table 2 (continued)

| Examples of selections form plasmid uptake and maintenance in <u>E. coli</u> |                                  |  |
|--|----------------------------------|--|
| gene   | (alternate designation)          | function   |
| Cam <sup>R</sup><br>colicin immunity<br>TrpA <sup>+</sup>                    | (Cmr <sup>R</sup> , <u>cat</u> ) | acetyltransferase<br>binds to colicin <u>in vivo</u><br>complementation of <u>trpA</u> |

Table 3

| Examples of selections for plasmid uptake and maintenance in <u>S. cerevisiae</u> |                            |
|---|----------------------------|
| gene  | function                   |
| Ura3 <sup>+</sup>   | complements ura3 auxotroph |
| Trp1 <sup>+</sup>   | complements trp1 auxotroph |
| Leu2 <sup>+</sup>   | complements leu2 auxotroph |
| His3 <sup>+</sup>   | complements his3 auxotroph |
| Neo <sup>R</sup>  | resistance to G418         |

Table 4: Agents for Selection  
of DBP Binding in E. coli and Relevant Genotypes

| Plasmid<br>Genotype  | Forward Selection |   | Reverse Selection             |  |
|--|-------------------|---|-------------------------------|--|
|  | Agent             | Host<br>Genotype  | Agent                         | Host<br>Genotype   |
| Galactose-1-phosphate uridylyltransferase and galactokinase                  |                   |   |                               |  |
| <u>galT</u> <sup>+</sup> &<br><u>galK</u> <sup>+</sup>                       | galactose         | <u>galE</u> <sup>-</sup> , <u>galT</u> <sup>-</sup> ,<br><u>galK</u> <sup>-</sup> | galactose as<br>sole C source | <u>galE</u> <sup>+</sup> &<br>( <u>galT</u> <sup>-</sup> or <u>galK</u> <sup>-</sup> ) |
| Tetracycline resistance ( <u>E. coli</u> K-12 strains are Tet <sup>S</sup> ) |                   |   |                               |  |
| <u>tetA</u> <sup>+</sup>   | fusaric acid      | Tet <sup>S</sup>  | Tc                            | Tet <sup>S</sup>   |
| beta galactosidase   |                   |   |                               |  |
| <u>lacZ</u> <sup>+</sup>   | phenylgalactoside | <u>lacZ</u> <sup>-</sup>  | lactose as<br>sole C source   | <u>lacZ</u> <sup>-</sup>   |



Table 4 (continued): Agents for Selection  
of DBP Binding in E. coli and Relevant Genotypes

| Plasmid<br>Genotype                          | Forward Selection |  | Reverse Selection   |  |
|--|-------------------|--|---|--|
|  | Agent             | Host<br>Genotype                               | Agent   | Host<br>Genotype   |
| Phe tRNA synthetase                          |                   |  |   |  |
| <u>pheS</u> <sup>+</sup>                     |                   | fluorophenylalanine <u>pheS12</u> <sup>+</sup> | growth at high<br>temperature                                       | <u>pheS</u> -amber,<br>sup-ts                                |
| Transport of arginine, lysine, and ornithine |                   |  |   |  |
| <u>argP</u> <sup>+</sup>                     | canavanine        | <u>argP</u> <sup>-</sup><br>Arg prototroph     | requirement for<br>arginine and<br>lysine at low<br>conc. in medium | <u>argP</u> <sup>-</sup> &<br>Arg auxotroph<br>Lys auxotroph |

Table 4 (continued): Agents for Selection  
of DBP Binding in *E. coli* and Relevant Genotypes

| Plasmid<br>Genotype                  | Forward Selection             |                          | Reverse Selection                                       |                          |
|--------------------------------------|-------------------------------|--------------------------|---|--------------------------|
|                                      | Agent                         | Host<br>Genotype         | Agent   | Host<br>Genotype         |
| thymidylate synthetase               |                               |                          |   |                          |
| <u>thyA</u> <sup>+</sup>             | trimethoprim +<br>thymidylate | <u>thyA</u> <sup>-</sup> | thymidylate<br>omitted from<br>defined medium           | <u>thyA</u> <sup>-</sup> |
| cAMP Receptor Protein (note 3)       |                               |                          |   |                          |
| <u>crp</u> <sup>+</sup>              | fosfomycin                    | <u>crp</u> <sup>-</sup>  | lactose or other<br>regulated sugar<br>as sole C source | <u>crp</u> <sup>-</sup>  |
| Orotidine-5'-phosphate decarboxylase |                               |                          |   |                          |
| <u>pyrF</u> <sup>+</sup>             | 5-fluoroorotate               | <u>pyrF</u> <sup>-</sup> | Thymine & cytosine<br>requirement on                    | <u>pyrF</u> <sup>-</sup> |

Table 4 (continued): Agents for Selection  
of DBP Binding in *E. coli* and Relevant Genotypes

| Plasmid<br>Genotype  | Forward Selection             |  | Reverse Selection           |  |
|--|-------------------------------|--|-----------------------------|--|
|  | Agent                         | Host<br>Genotype                                       | Agent                       | Host<br>Genotype                                       |
| mannosephosphotransferase enzyme II  |                               |  |                             |  |
| <u>ptsM</u> <sup>+</sup>   | deoxyglucose                  | <u>ptsM</u> <sup>-</sup>                               | Mannose as<br>sole C source | <u>ptsM</u> <sup>-</sup>                               |
| Fusion protein   |                               |  |                             |  |
| <u>secA</u> <sup>+</sup> &<br><u>malE</u><br>signal-<br><u>lacZ</u> fusion | lactose as sole C<br>(note 2) | <u>secA</u> <sup>-</sup> &<br><u>lacZ</u> <sup>-</sup> | phenylgalactoside           | <u>secA</u> <sup>-</sup> &<br><u>lacZ</u> <sup>-</sup> |

Table 4 (continued): Agents for Selection  
of DBP Binding in E. coli and Relevant Genotypes

| Plasmid<br>Genotype             | Forward Selection |                          | Reverse Selection   |  |
|---------------------------------|-------------------|--------------------------|---|--|
|                                 | Agent             | Host<br>Genotype         | Agent   | Host<br>Genotype   |
| Outer membrane protein (note 4) |                   |                          |   |  |
| <u>ompA</u> <sup>+</sup>        | colicin E1        | <u>ompA</u> <sup>-</sup> | HfrH( <u>thr</u> <sup>+</sup> , <u>leu</u> <sup>+</sup> ,<br><u>str</u> <sup>S</sup> )<br>conjugation | <u>thr</u> <sup>-</sup> , <u>leu</u> <sup>-</sup> ,<br><u>ompA</u> , <u>str</u> <sup>R</sup> |
|                                 | colicin E2        |                          |   |  |
|                                 | colicin E3        |                          |   |  |
|                                 | phage TuII        |                          |   |  |
|                                 | phage K3          |                          |   |  |
|                                 | phage 4-59        |                          |   |  |

Table 4 (continued): Agents for Selection  
of DBP Binding in E. coli and Relevant Genotypes

| Plasmid<br>Genotype      | Forward Selection |  | Reverse Selection                           |   |
|--------------------------|-------------------|--|---|---|
|                          | Agent             | Host<br>Genotype                           | Agent                                       | Host<br>Genotype                            |
| Vitamin B12 transport    |                   |  |   |   |
| <u>btuB</u> <sup>+</sup> | phage BF23        | <u>btuB</u> <sup>-</sup><br>B12 prototroph | requirement for<br>B12 in defined<br>medium | <u>btuB</u> <sup>-</sup> &<br>B12 auxotroph |

Table 4 (continued): Agents for Selection  
of DBP Binding in *E. coli* and Relevant Genotypes

| Plasmid<br>Genotype      | Forward Selection |                          | Reverse Selection                                      |                          |
|--------------------------|-------------------|--------------------------|--|--------------------------|
|                          | Agent             | Host<br>Genotype         | Agent  | Host<br>Genotype         |
| Maltose transport        |                   |                          |  |                          |
| <u>lamB</u> <sup>+</sup> | Phage $\lambda$   | <u>lamB</u> <sup>-</sup> | growth on<br>maltose as sole<br>C source               | <u>lamB</u> <sup>-</sup> |
| Ferrichrome receptor     |                   |                          |  |                          |
| <u>tonA</u> <sup>+</sup> | Phage phi80       | <u>tonA</u> <sup>-</sup> | Requirement for<br>Fe hydroxamate as<br>sole Fe source | <u>tonA</u> <sup>-</sup> |
| Colicin I receptor       |                   |                          |  |                          |
| <u>cir</u> <sup>+</sup>  | colicin I         | <u>cir</u> <sup>-</sup>  | screen for<br>colicin I resistance                     | <u>cir</u> <sup>-</sup>  |

Table 4 (continued): Agents for Selection  
of DBP Binding in *E. coli* and Relevant Genotypes

| Plasmid<br>Genotype   | Forward Selection                        |                          | Reverse Selection                                   |  |
|---|--|--------------------------|---|--|
|   | Agent                                    | Host<br>Genotype         | Agent   | Host<br>Genotype   |
| Nucleoside uptake, colicin K receptor, phage T6 receptor (note 5) |  |                          |   |  |
| <u>tsx</u> <sup>+</sup>   | colicin K                                | <u>tsx</u> <sup>-</sup>  | requirement for<br>nucleosides in<br>defined medium | <u>tsx</u> <sup>-</sup><br>thymine auxotroph<br>purine auxotroph |
|   | Phage T6                                 |                          |   |  |
| Aromatic amino acid transport                                     |  |                          |   |  |
| <u>arop</u> <sup>+</sup>  | thienylalanine or<br>fluorophenylalanine | <u>arop</u> <sup>-</sup> | requirement for<br>tryptophan in<br>defined medium  | Trp auxotroph<br><u>arop</u> <sup>-</sup>                        |
| Cysteine synthetase   |  |                          |   |  |
| <u>cysK</u> <sup>+</sup>  | selenate or<br>azaserine in medium       | <u>cysK</u> <sup>-</sup> | growth on medium<br>lacking cysteine                | <u>cysK</u> <sup>-</sup>   |

Table 4 (continued): Agents for Selection  
of DBP Binding in *E. coli* and Relevant Genotypes

| Plasmid<br>Genotype            | Forward Selection |                          | Reverse Selection  |                          |
|--------------------------------|-------------------|--------------------------|--|--------------------------|
|                                | Agent             | Host<br>Genotype         | Agent  | Host<br>Genotype         |
| C4 dicarboxylic acid transport |                   |                          |  |                          |
| <u>dctA</u> <sup>+</sup>       | 3-fluoromalate    | <u>dctA</u> <sup>-</sup> | grow on C4<br>dicarboxylic<br>acids as C and energy<br>source. | <u>dctA</u> <sup>-</sup> |
| Spectinomycin                  |                   |                          |  |                          |
| <u>aadA</u> occluded           | Spectinomycin     | any                      | none   |                          |



Table 4 (continued): Agents for Selection  
of DBP Binding in E. coli and Relevant Genotypes

Notes:

1) Deletions are strongly preferred over point mutations.

2) Only secA gene need be controlled by DBP.

3) Mutations in crp are highly pleotropic; some effects seen in cell wall. crp best used in connection with selections having intracellular action.

4) Resistance to colicins can arise in several ways. Use of two or more E-colicins discriminates against other mechanisms. Because colicins do not replicate they are preferred over phage for selection. Phages are useful to verify selection of cells repressing expression of ompA.

5) Because colicins do not replicate, they are preferred over phage for selection. Phages are used to verify selection of cells repressing expression of tsx.

Table 5: Some Recommended Pairs of Selectable  
Binding Marker Genes

## A) Recommended pairs:

|                         |             |  |                    |
|-------------------------|-------------|--|--------------------|
| <u>galT.K</u>           | <u>tetA</u> |  | <u>galT.KpheS</u>  |
| <u>argP</u>             | <u>pheS</u> |  | <u>tetA thyA</u>   |
| <u>lacZ</u>             | <u>tetA</u> |  | <u>ptsM thyA</u>   |
| <u>dctA</u>             | <u>cysK</u> |  | <u>ompA pyrF</u>   |
| <u>crp</u>              | <u>thyA</u> |  | <u>btuB pyrF</u>   |
| <u>lamB</u>             | <u>thyA</u> |  | <u>tonA galT.K</u> |
| <u>secA</u> &           | <u>pyrF</u> |  | <u>cir cysK</u>    |
| <u>malE-lacZ</u> fusion |             |  |                    |
| <u>tsx</u>              | <u>cysK</u> |  | <u>aroP lacZ</u>   |
| <u>dctA</u>             | <u>thyA</u> |  |                    |

## B) Less Preferred pairs:

|                         |               |
|-------------------------|---------------|
| <u>tetA</u>             | <u>argP</u>   |
| <u>secA</u> &           | <u>lacZ</u>   |
| <u>malE-lacZ</u> fusion |               |
| <u>pyrF</u>             | <u>thyA</u>   |
| <u>lamB</u>             | <u>galT.K</u> |
| <u>cir</u>              | <u>tsx</u>    |
| <u>ptsM</u>             | <u>tetA</u>   |
| <u>tonA</u>             | <u>ptsM</u>   |
| <u>crp</u>              | <u>lacZ</u>   |

## Reason

Both transport related.

Both related to lacZ  
function.

Both related to thymine

Both related to sugar me  
bolism.

Both related to colicin

Both transport related

Both transport related

Both related to sugar me  
bolism

Table 6 Promot rs

A: Correlation between Sequence Homology and Promoter Strength (MULL84)

| Promoter              | Homology score | Log K <sub>pk2</sub> |
|-----------------------|----------------|----------------------|
| T7 A1                 | 74.0           | 7.40                 |
| T7 A2                 | 73.4           | 7.20                 |
| λ P <sub>R</sub>      | 58.6           | 7.13                 |
| lac UV5               | 59.2           | 6.94                 |
|                       | 59.2           | 6.30                 |
| T7 D                  | 63.9           | 6.30                 |
|                       | 63.9           | 6.00                 |
| Tn10 P <sub>out</sub> | 56.2           | 6.71                 |
| Tn10 P <sub>in</sub>  | 52.1           | 6.18                 |
| λ P <sub>RM</sub>     | 49.7           | 4.71                 |
|                       | 49.7           | 4.17                 |
| P <sub>amp</sub>      | 52.7           |                      |
| P <sub>neo</sub>      | 58.0           |                      |

Table 6 (continued) Promoters

## B: Sequences of some promoters

| Name                      | -35            | -10                    | +1    |
|---------------------------|----------------|------------------------|-------|
| T7 A1                     | GTATTGACTTAAAG | TCTAACCTATAGGATACTTAC  | AGCCA |
| T7 A2                     | GTATTGACAACATG | AAGTAACATGCAGTAAGATACA | AATCG |
| $\lambda$ P <sub>R</sub>  | GTGTTGACTATTTT | ACCTCTGGCGGTAGAATGGT   | TGCA  |
| lac UV5                   | GGCTTTACACTTTA | TGCTTCCGGCTCATATAATGTG | TGGA  |
| T7 D                      | GCGTTGACTTGATG | GGTCTTTATGTGTAGGCTTTA  | GGTG  |
| Tn10 P <sub>out</sub>     | GGGCAGAATTGGTA | AAGAGAGTCGTGTAAAATATC  | GAGT  |
| Tn10 P <sub>in</sub>      | AGGTGGATACACAT | CTTGTCATATGATCAAATGGT  | TTCG  |
| $\lambda$ P <sub>RM</sub> | TGTTAGATATTTAT | CCCTTGCGGTGATAGATTTAA  | CATA  |
| P <sub>amp</sub>          | ACATTCAAATATGT | ATCCGCTCATGAGACAATAAC  | CCTG  |
| P <sub>neo</sub>          | GAATTGCCAGCTGG | GGCGCCCTCTGGTAAGGTGG   | GAAG  |

Table 7      FUNCTIONAL SUBSTITUTIONS IN HELIX 5 OF  $\lambda$   
REPRESSOR

5

|    | 84          | 85     | 86     | 87     | 88     | 89     | 90     | 91     |
|----|-------------|--------|--------|--------|--------|--------|--------|--------|
|    | .           | .      | .      | .      | .      | .      | .      | .      |
| 10 | -----I----- | Y----- | E----- | M----- | Y----- | E----- | A----- | V----- |
|    | .           | .      | .      | .      | .      | .      | .      | .      |
|    | I           | I      | I      | M      | I      | I      | V      | I      |
| 15 |             | V      | F      | L      | F      | L      | L      | V      |
|    |             | L      | L      |        | V      | M      | W      | L      |
|    |             | W      | W      |        | W      | A      | M      | A      |
| 20 |             | A      | M      |        | M      | G      | A      | C      |
|    |             | G      | A      |        | A      | T      | G      | S      |
|    |             | C      | G      |        | C      | S      | S      | T      |
|    |             | Y      | Y      |        | Y      | E      | H      |        |
| 25 |             | T      | T      |        | S      | Q      | Q      |        |
|    |             | S      | S      |        |        | D      |        |        |
|    |             | E      | E      |        |        | K      |        |        |
| 30 |             | Q      | Q      |        |        | R      |        |        |
|    |             | R      | D      |        |        |        |        |        |

35

Table 8  
Some Preferred Initial DBPs

40

|    |   |
|----|---|
|    | $\lambda$ cI repressor                    |
|    | $\lambda$ Cro                             |
|    | 434 cI repressor                          |
| 45 | 434 Cro                                   |
|    | P22 Mnt                                   |
|    | P22 Arc                                   |
|    | P22 cII repressor                         |
| 50 | $\lambda$ cII repressor                   |
|    | $\lambda$ Xis                             |
|    | $\lambda$ Int                             |
| 55 | cAMP Receptor Protein from <u>E. coli</u> |
|    | Trp Repressor from <u>E. coli</u>         |
|    | Kr protein from <u>Drosophila</u>         |

Table 8 (continued)

Some Preferred Initial DBPs

---

Transcription Factor IIIA from Xenopus laevis

5

Lac Repressor from E. coli

Tet Repressor from Tn10

Mu repressor from phage mu

10

Yeast MAT-a1-alpha2

Polyoma Large T antigen

SV40 Large T antigen

Adenovirus E1A

15

Human Transcription Factor SP1 (a zinc finger protein)

Human Transcription Factor AP1 (product of jun)

[0405] Table 9, Table 10, and Table 11 have been deleted.

20

25

30

35

40

45

50

55

**Table 12 MISSENSE MUTATIONS IN  $\lambda$  REPRESSOR AMINO-TERMINAL DOMAIN**

5 (S)  
(Q)  
(Y)  
(L)

10 Q 5 (W) 1 (K) 1 2 3  
0 P | 5 0 5 0  
|-----|  
S T K K P L T O E O L E D A R R L K A I Y E K K K N E L G  
alpha 1 | | F (F) P  
E (H) D C S

15





5  
10  
15  
20  
25  
30  
35  
40  
45  
50  
55

|   |   |   |     |         |     |       |   |  |
|---|---|---|-----|---------|-----|-------|---|--|
| 5 | 6   | 7 | 7   | 8       | K*  | (Q)   | 9 |  |
|   | 5   | 0 | 5   | 0       | (W) | (H(S) | 0 |  |
|   | ----- .   |   |     | -----   |     |       |   |  |
|   | N A L L A K I L K V S V E E F S P S I A R E I Y E M Y E A V S |   |     |         |     |       |   |  |
|   | alpha 4   |   |     | alpha 5 |     |       |   |  |
|   | R   |   | N   | S       |     |       |   |  |
|   | (C)   |   | R   |         |     |       |   |  |
|   |   |   | I   |         |     |       |   |  |
|   |   |   | (T) |         |     |       |   |  |
|   |   |   | (C) |         |     |       |   |  |

Substitutions occurring at solvent exposed positions in the unbound repressor dimer are shown above the wild type sequence.

Substitutions occurring at internal positions are shown below the wild type sequence.

Substitutions that produce repressor dimers with normal or nearly normal DNA binding affinities are shown in parentheses.

Substitutions that produce repressor dimers with increased DNA binding affinities are designated by \*.

Table 13. MISSENSE MUTATIONS IN P22 ARC REPRESSOR  
THAT PRODUCE AN ARC<sup>-</sup> PHENOTYPE

|    |                     |    |    |    |    |    |    |   |   |   |   |   |    |   |   |   |   |   |   |   |    |   |   |   |   |   |   |    |     |   |    |
|----|---------------------|----|----|----|----|----|----|---|---|---|---|---|----|---|---|---|---|---|---|---|----|---|---|---|---|---|---|----|-----|---|----|
| 5  |                     | 5  | 10 | 15 | 20 | 25 | 30 |   |   |   |   |   |    |   |   |   |   |   |   |   |    |   |   |   |   |   |   |    |     |   |    |
|    | .                   | .  | .  | .  | .  | .  | .  |   |   |   |   |   |    |   |   |   |   |   |   |   |    |   |   |   |   |   |   |    |     |   |    |
| 10 | M                   | K  | G  | M  | S  | K  | M  | P | Q | F | N | L | R  | W | P | R | E | V | L | D | L  | V | R | K | V | A | E | E  | N   | G |    |
|    | <u>high yield</u>   |    |    |    |    |    |    |   |   |   |   |   |    |   |   |   |   |   |   |   |    |   |   |   |   |   |   |    |     |   |    |
|    |                     | Q  | R  | I  | C  |    |    |   |   |   |   |   | L  |   |   |   |   |   |   |   |    |   |   |   |   |   |   |    |     |   |    |
| 15 |                     | T  | K  |    |    |    |    |   |   |   |   |   | V  |   |   |   |   |   |   |   |    |   |   |   |   |   |   |    |     |   |    |
|    | <u>medium yield</u> |    |    |    |    |    |    |   |   |   |   |   |    |   |   |   |   |   |   |   |    |   |   |   |   |   |   |    |     |   |    |
|    |                     | R  |    |    |    |    |    |   |   |   |   |   |    |   |   |   | A | G |   |   | A  | F |   |   | G |   |   | A  |     |   |    |
| 20 |                     |    |    |    |    |    |    |   |   |   |   |   |    |   |   |   | E |   |   |   |    |   |   |   |   |   |   |    |     |   |    |
|    | <u>low yield</u>    |    |    |    |    |    |    |   |   |   |   |   |    |   |   |   |   |   |   |   |    |   |   |   |   |   |   |    |     |   |    |
|    |                     | R  |    |    |    |    |    |   |   |   |   |   |    |   |   |   |   |   |   |   | N  |   | H |   | I | T |   | K  | K   |   |    |
| 25 |                     |    |    |    |    |    |    |   |   |   |   |   |    |   |   |   |   |   |   |   | Y  |   | C |   | V |   |   | T  |     |   |    |
|    |                     |    |    |    |    |    |    |   |   |   |   |   |    |   |   |   |   |   |   |   |    |   |   |   |   |   |   |    | Y   |   |    |
|    | <u>undetermined</u> |    |    |    |    |    |    |   |   |   |   |   |    |   |   |   |   |   |   |   |    |   |   |   |   |   |   |    |     |   |    |
| 30 |                     | L  |    |    |    |    |    |   |   |   |   |   |    |   |   |   |   |   |   |   | G  |   | S |   | V |   |   |    |     |   |    |
|    |                     |    |    |    |    |    |    |   |   |   |   |   |    |   |   |   |   |   |   |   |    |   |   |   |   |   |   |    |     |   |    |
| 35 |                     | 35 |    |    |    |    |    |   |   |   |   |   | 40 |   |   |   |   |   |   |   | 45 |   |   |   |   |   |   | 50 |     |   | 55 |
|    | .                   | .  | .  | .  | .  | .  | .  | . | . | . | . | . | .  | . | . | . | . | . | . | . | .  | . | . | . | . | . | . | .  | .   | . | .  |
|    | R                   | S  | V  | N  | S  | E  | I  | Y | Q | R | V | M | E  | S | F | K | K | E | G | R | I  | G | A | - | - |   |   |    |     |   |    |
| 40 | <u>high yield</u>   |    |    |    |    |    |    |   |   |   |   |   |    |   |   |   |   |   |   |   |    |   |   |   |   |   |   |    |     |   |    |
|    |                     |    |    |    |    |    |    |   |   |   |   |   |    |   |   |   |   |   |   |   |    |   |   |   |   |   |   |    | Y+S |   |    |
|    | <u>medium yield</u> |    |    |    |    |    |    |   |   |   |   |   |    |   |   |   |   |   |   |   |    |   |   |   |   |   |   |    |     |   |    |
| 45 |                     | A  |    |    |    |    |    |   |   |   |   |   |    |   |   |   |   |   |   |   | A  |   | T |   |   |   |   |    |     |   |    |
|    | <u>low yield</u>    |    |    |    |    |    |    |   |   |   |   |   |    |   |   |   |   |   |   |   |    |   |   |   |   |   |   |    |     |   |    |
|    | W                   | F  | G  | H  |    | A  | M  | S | P | Q | A |   | G  |   | S |   |   |   |   |   | K  |   | P |   |   |   |   |    |     |   |    |
| 50 | L                   |    | K  |    | K  |    | D  |   |   |   |   |   |    |   | C |   |   |   |   |   |    |   |   |   |   |   |   |    |     |   |    |
|    |                     |    |    |    | G  |    | C  |   |   |   |   |   |    |   |   |   |   |   |   |   |    |   |   |   |   |   |   |    |     |   |    |

TABLE 14

MISSENSE MUTATIONS AT SOLVENT EXPOSED POSITIONS  
OF THE H-T-H REGIONS OF REPRESSOR PROTEINS

Table 14a\ Repressor

|  | (S)                   |                           | (N)  |                |    |
|--|-----------------------|---------------------------|------|----------------|----|
|  | L                     |                           | R    |                | C  |
|  | Y(K)T                 |                           |      |                |    |
|  | 35                    | 40                        | 45   | 50             | 55 |
|  | .                     | .                         | .    | .              | .  |
|  | <u>HELIX 2</u>        |                           | TURN | <u>HELIX 3</u> |    |
|  | Q E S V A D K M G M G | Q S G V G A L F N G I N A |      |                |    |
|  | *                     | *                         | *    | *              | *  |
|  | S P                   | E Y L                     | D D  | D D            | K  |
|  | K                     | L                         | V    |                |    |
|  | L                     | S                         |      |                |    |

Table 14b \ Cro

|  | F                       |                           | K    |                |    |
|--|-------------------------|---------------------------|------|----------------|----|
|  | K                       |                           | R T  |                |    |
|  | 15                      | 20                        | 25   | 30             | 35 |
|  | .                       | .                         | .    | .              | .  |
|  | <u>HELIX 2</u>          |                           | TURN | <u>HELIX 3</u> |    |
|  | G Q T K T A K D L G V Y | Q S A I N K A I H A G R K |      |                |    |
|  | *                       | *                         | *    | *              | *  |
|  | R H                     | D H N                     | N    | Q T            |    |
|  | E P                     | N L R                     | T    | L              |    |
|  |                         | C A                       | Q    |                |    |

Table 14C 434 Repressor

```

      5                                     H
                                           L
                                           V
                                           T
                                           A
      10
      20          25          30          35
      .           .           .           .
HELIX 2       TURN       HELIX 3
Q A E L A Q K V G T T Q Q S I E Q L E N
★                ★ ★              ★
      20
                                A
                                H
                                L
                                S
                                M
                                R
                                P
      30
                                K

```

Table 14d Trp Repressor

[illegible]

## Table 14, NOTES:

5 Positions in wild type repressors believed to contact DNA are indicated by a \* below the wild type residue.

10 Substitutions that greatly decrease repressor binding to DNA are shown below the wild type sequence.

Substitutions that produce repressors with normal or nearly normal DNA binding affinities are shown below the wild type sequence.

15 Substitutions that increase repressor affinity for DNA are shown in parentheses above the wild type sequence.

20  
Table 15: deleted.

Table 16: Genetic Code Table  
With Secondary-Structure  
Pr eferences

|       |  | Second Base |     |   |     |      |     |      |     |       |
|-------|--|-------------|-----|---|-----|------|-----|------|-----|-------|
| First |  | T           |     | C |     | A    |     | G    |     | Third |
| Base  |  |             |     |   |     |      |     |      |     | base  |
| T     |  | F           | b/a | S | a/b | Y    | b   | C    | b   | T     |
|       |  | F           | b/a | S | a/b | Y    | b   | C    | b   | C     |
|       |  | L           | a/b | S | a/b | stop |     | stop |     | A     |
|       |  | L           | a/b | S | a/b | stop |     | W    | b/a | G     |
| C     |  | L           | a/b | P | -   | H    | a/b | R    | b/a | T     |
|       |  | L           | a/b | P | -   | H    | a/b | R    | b/a | C     |
|       |  | L           | a/b | P | -   | Q    | b/a | R    | b/a | A     |
|       |  | L           | a/b | P | -   | Q    | b/a | R    | b/a | G     |
| A     |  | I           | b   | T | b   | N    | a/b | S    | a/b | T     |
|       |  | I           | b   | T | b   | N    | a/b | S    | a/b | C     |
|       |  | I           | b   | T | b   | K    | a/b | R    | b/a | A     |
|       |  | M           | b   | T | b   | K    | a/b | R    | b/a | G     |
| G     |  | V           | b   | A | a   | D    | a/b | G    | b/a | T     |
|       |  | V           | b   | A | a   | D    | a/b | G    | b/a | C     |
|       |  | V           | b   | A | a   | E    | a   | G    | b/a | A     |
|       |  | V           | b   | A | a   | E    | a   | G    | b/a | G     |

Amino acids denoted "b" strongly favor extended structures.  
Amino acids denoted "b/a" favor extended structures.  
Amino acids denoted "a/b" strongly favor helical structures.  
Amino acids denoted "a" very strongly favor helices.  
Proline is denoted "-" and favors neither beta sheets nor helices.

b: I, M, V, T, Y, C  
b/a: F, Q, R, G, W  
a/b: L, S, H, N, K, D  
a: A, E  
-: P

Table 17

| Fraction of DNA molecules having n non-parental bases when reagents<br>that have fraction M of parental nucleotide.                           |         |                    |        |                    |        |          |
|---|---------|--------------------|--------|--------------------|--------|----------|
| Number of bases using mixed reagents is 30.   |         |                    |        |                    |        |          |
| M   | .9965   | .97716             | .92612 | .8577              | .79433 | .63096   |
| f0  | .9000   | .5000              | .1000  | .0100              | .0010  | .000001  |
| f1  | .09499  | .35061             | .2393  | .04977             | .00777 | .0000175 |
| f2  | .00485  | .1188              | .2768  | .1197              | .0292  | .000149  |
| f3  | .00016  | .0259              | .2061  | .1854              | .0705  | .000812  |
| f4  | .000004 | .00409             | .1110  | .2077              | .1232  | .003207  |
| f8  | 0.      | $2 \times 10^{-7}$ | .00096 | .0336              | .1182  | .080165  |
| f16   | 0.      | 0.                 | 0.     | $5 \times 10^{-7}$ | .00006 | .027281  |
| f23   | 0.      | 0.                 | 0.     | 0.                 | 0.     | .0000089 |
| most  | 0       | 0                  | 2      | 5                  | 7      | 12       |
| fn is the fraction of all synthetic DNA molecules having n non-parental<br>bases.<br>"most" is the value of n having the highest probability. |         |                    |        |                    |        |          |

Table 18: best vgCodon

```

5      Program "Find Optimum vgCod n."
      INITIALIZE-MEMORY-OF-ABUNDANCES
      DO ( t1 = 0.21 to 0.31 in steps of 0.01 )
10      . DO ( c1 = 0.13 to 0.23 in steps of 0.01 )
      . . DO ( a1 = 0.23 to 0.33 in steps of 0.01 )
      Comment      calculate g1 from other concentrations
      . . . g1 = 1.0 - t1 - c1 - a1
15      . . . IF( g1 .ge. 0.15 )
      . . . . DO ( a2 = 0.37 to 0.50 in steps of 0.01 )
      . . . . . DO ( c2 = 0.12 to 0.20 in steps of 0.01 )
20      Comment      Force D+E = R + K
      . . . . . g2 = (g1*a2 -.5*a1*a2)/(c1+0.5*a1)
      Comment      Calc t2 from other concentrations.
25      . . . . . t2 = 1. - a2 - c2 - g2
      . . . . . IF(g2.gt. 0.1.and. t2.gt.0.1)
      . . . . . CALCULATE-ABUNDANCES
      . . . . . COMPARE-ABUNDANCES-TO-PREVIOUS-ONES
30      . . . . . ..end_IF_block
      . . . . . ..end_DO_loop ! c2
      . . . . . ..end_DO_loop ! a2
35      . . . ..end_IF_block ! if g1 big enough
      . . . ..end_DO_loop ! a1
      . . . ..end_DO_loop ! c1
40      ..end_DO_loop ! t1
      WRITE the best distribution and the abundances.

```



Table 19: Abundances obtained  
from optimum vgCodon

5

| Amino acid | Abundance  | Amino acid | Abundance  |
|------------|------------|------------|------------|
| A          | 4.80%      | C          | 2.86%      |
| D          | 6.00%      | E          | 6.00%      |
| F          | 2.86%      | G          | 6.60%      |
| H          | 3.60%      | I          | 2.86%      |
| K          | 5.20%      | L          | 6.82%      |
| M          | 2.86%      | N          | 5.20%      |
| P          | 2.88%      | Q          | 3.60%      |
| R          | 6.82%      | S          | 7.02% mfaa |
| T          | 4.16%      | V          | 6.60%      |
| W          | 2.86% lfaa | Y          | 5.20%      |
| stop       | 5.20%      |            |            |

25

lfaa = least-favored amino acid  
mfaa = most-favored amino acid  
ratio =  $\text{Abun}(W)/\text{Abun}(S) = 0.4074$

30

35

| i | $(1/\text{ratio})^j$ | $(\text{ratio})^j$    | stop-free |
|---|----------------------|-----------------------|-----------|
| 1 | 2.454                | .4074                 | .9480     |
| 2 | 6.025                | .1660                 | .8987     |
| 3 | 14.788               | .0676                 | .8520     |
| 4 | 36.298               | .0275                 | .8077     |
| 5 | 89.095               | .0112                 | .7657     |
| 6 | 218.7                | $4.57 \times 10^{-3}$ | .7258     |
| 7 | 536.8                | $1.86 \times 10^{-3}$ | .6881     |

40

45

50

55

Table 20: Calculate worst codon.

Program "Find worst vgCodon within Serr of given distribution."

INITIALIZE-MEMORY-OF-ABUNDANCES

READ Serr                      Comment Serr is % error level.

Comment T1i,C1i,A1i,G1i; T2i,C2i,A2i,G2i, T3i,G3i

Comment are the intended nt-distribution.

READ T1i, C1i, A1i, G1i

READ T2i, C2i, A2i, G2i

READ T3i, G3i

Fdwn = 1.-Serr

Fup = 1.+Serr

DO ( t1 = T1i\*Fdwn to T1i\*Fup in 7 steps)

. DO ( c1 = C1i\*Fdwn to C1i\*Fup in 7 steps)

. . DO ( a1 = A1i\*Fdwn to A1i\*Fup in 7 steps)

. . . g1 = 1. - t1 - c1 - a1

. . . IF( (g1-G1i)/G1i .lt. -Serr)

Comment g1 too far below G1i, push it back

. . . . g1 = G1i\*Fdwn

. . . . factor = (1.-g1)/(t1 + c1 + a1)

. . . . t1 = t1\*factor

. . . . c1 = c1\*factor

. . . . a1 = a1\*factor

. . . . .end\_IF\_block

. . . IF( (g1-G1i)/G1i .gt. Serr)

Comment g1 too far above G1i, push it back

. . . . g1 = G1i\*Fup

. . . . factor = (1.-g1)/(t1 + c1 + a1)

. . . . t1 = t1\*factor

. . . . c1 = c1\*factor

. . . . a1 = a1\*factor

. . . . .end\_IF\_block

. . . DO ( a2 = A2i\*Fdwn to A2i\*Fup in 7 steps)

. . . . DO ( c2 = C2i\*Fdwn to C2i\*Fup in 7 steps)

. . . . . DO (g2=G2i\*Fdwn to G2i\*Fup in 7 steps)

Comment Calc t2 from other concentrations.

. . . . . t2 = 1. - a2 - c2 - g2

. . . . . IF( (t2-T2i)/T2i .lt. -Serr)

Table 20, continued: Calculate worst codon.

Comment t2 too far below T2i, push it back

```

5      . . . . . t2 = T2i*Fdwn
      . . . . . factor = (1.-t2)/(a2 + c2 + g2)
      . . . . . a2 = a2*factor
      . . . . . c2 = c2*factor
10     . . . . . g2 = g2*factor
      . . . . . ..end_IF_block
      . . . . . IF( (t2-T2i)/T2i .gt. Serr)

```

15 Comment t2 too far above T2i, push it back

```

      . . . . . t2 = T2i*Fup
      . . . . . factor = (1.-t2)/(a2 + c2 + g2)
      . . . . . a2 = a2*factor
20     . . . . . c2 = c2*factor
      . . . . . g2 = g2*factor
      . . . . . ..end_IF_block
      . . . . . IF(g2.gt. 0.0 .and. t2.gt.0.0)
25     . . . . . t3 = 0.5*(1.-Serr)
      . . . . . g3 = 1. - t3
      . . . . . CALCULATE-ABUNDANCES
      . . . . . COMPARE-ABUNDANCES-TO-PREVIOUS-ONES
30     . . . . . t3 = 0.5
      . . . . . g3 = 1. - t3
      . . . . . CALCULATE-ABUNDANCES
35     . . . . . COMPARE-ABUNDANCES-TO-PREVIOUS-ONES
      . . . . . t3 = 0.5*(1.+Serr)
      . . . . . g3 = 1. - t3
40     . . . . . CALCULATE-ABUNDANCES
      . . . . . COMPARE-ABUNDANCES-TO-PREVIOUS-ONES
      . . . . . ..end_IF_block
      . . . . . ..end_DO_loop ! g2
45     . . . . . ..end_DO_loop ! c2
      . . . . . ..end_DO_loop ! a2
      . . . . . ..end_DO_loop ! a1
50     . . . . . ..end_DO_loop ! c1
      . . . . . ..end_DO_loop ! t1
      WRITE the WORST distribution and the abundances.

```

55

Table 21: Abundances obtained  
using optimum vgCodon assuming  
5% errors

| Amino acid | Abundance  | Amino acid | Abundance |
|------------|------------|------------|-----------|
| A          | 4.59%      | C          | 2.76%     |
| D          | 5.45%      | E          | 6.02%     |
| F          | 2.49% lfaa | G          | 6.63%     |
| H          | 3.59%      | I          | 2.71%     |
| K          | 5.73%      | L          | 6.71%     |
| M          | 3.00%      | N          | 5.19%     |
| P          | 3.02%      | Q          | 3.97%     |
| R          | 7.68% mfaa | S          | 7.01%     |
| T          | 4.37%      | V          | 6.00%     |
| W          | 3.05%      | Y          | 4.77%     |
| stop       | 5.27%      |            |           |

$$\text{ratio} = \text{Abun}(F)/\text{Abun}(R) = 0.3248$$

| i | $(1/\text{ratio})^i$ | $(\text{ratio})^i$    | stop-free |
|---|----------------------|-----------------------|-----------|
| 1 | 3.079                | .3248                 | .9473     |
| 2 | 9.481                | .1055                 | .8973     |
| 3 | 29.193               | .03425                | .8500     |
| 4 | 89.888               | .01112                | .8052     |
| 5 | 276.78               | $3.61 \times 10^{-3}$ | .7627     |
| 6 | 852.22               | $1.17 \times 10^{-3}$ | .7225     |
| 7 | 2624.1               | $3.81 \times 10^{-4}$ | .6844     |

Table 22, deleted

## Tables for Example 1

Table 100:  $\lambda$  Or3 Downstream  
of Pamp that promotes galT,K

5' GAT|CGT|TAA|CGG|GCC|CTT|CTA|AAT|ACA|TTC|AAA|-  
 olig#4 3' ca att gcc cgg gaa gat tta tgt aag ttt|-  
HpaI || ApaI | -35 |

TAT|GTA|TCC|GCT|CAT|GAG|ACA|ATA|ACC|-  
ata cat agg cga gta ctc tgt tat tgg-  
-10 |

CTT|ATC|ACC|GCA|AGG|GAT|ATC|TAG|AGT|C 3' = olig#3  
gaa tag tgg cgt tcc cta tag atc t 5'  
 $\lambda$  Or3 || XbaI |

Table 101:  $\lambda$  O<sub>R</sub>3 Downstr am  
of P<sub>neo</sub> that promotes tet

5' | CCT|GCG|AAC|CGG|AAT|TGC|CAG|-  
Olig#6 3' | ggc cgc ttg gcc tta acg gtc-  
| StuI | | -35 |

| CTG|GGG|CGC|CCT|CTG|GTA|AGG|TTG|-  
gac ccc gcg gga gac cat tcc aac-  
| -10 |

| GGA|TAT|CAC|CGC|AAG|GGA|TA 3' = Olig#5  
ggg ata gtg gcg ttc cct att cg a 5'  
|  $\lambda$  O<sub>R</sub>3 | HindIII |

Table 102: ray gene  
using lacUV5 as promoter

SpeI-BstEII-(BalI-PpuMI-BglII-BamHI-AvaI)  
-KpnI-(Trp terminator)-SfiI; !

5'-ACTAGT CCAGG C TTTACA CTT TATGC TTCCG GCTCG TATAAT GTGT GG  
|SpeI|

AAT TGTGA GCGGA TAACA ATTTC ACAC ! lacUV5

A GGTAACC AGGAGGAAATAAA ! BstEII & Shine-Dalgarno seq.  
|BstE2|

|   |   |   |   |   |   |   |   |   |    |    |    |    |
|---|---|---|---|---|---|---|---|---|----|----|----|----|
| m | e | q | r | i | t | l | k | d | y  | a  | m  | r  |
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |

ATG GAA CAA CGC ATA ACC CTA AAG GAC TAC GCG ATG CGC

|    |    |    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|----|----|
| f  | g  | q  | t  | k  | t  | a  | k  | d  | l  |
| 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 |

TTT GGC CAA ACC AAG ACA GCG AAG GAC CTA  
|Bal I| |PpuM I|

|    |    |    |    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|----|----|----|
| g  | v  | y  | q  | s  | a  | i  | n  | k  | a  | i  |
| 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 | 32 | 33 | 34 |

GGG GTG TAT CAG AGC GCG ATT AAC AAG GCC ATC

|    |    |    |    |    |    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|----|----|----|----|----|
| h  | a  | g  | r  | k  | i  | f  | l  | t  | i  | n  | a  | d  |
| 35 | 36 | 37 | 38 | 39 | 40 | 41 | 42 | 43 | 44 | 45 | 46 | 47 |

CAT GCC GGC CGA AAG ATC TTC CTA ACC ATT AAC GCT GAT  
|Bgl II|

Table 102, continued

g s v y a e e v k p f p s  
 48 49 50 51 52 53 54 55 56 57 58 59 60  
 GGA TCC GTC TAC GCG GAA GAG GTA AAG CCC TTC CCG AGT  
BamHI Ava I

n k k t t a . . .  
 61 62 63 64 65 66 67 67 68  
 AAC AAA AAA ACA ACA GCG TAA TAG TA GGTACC  
KpnI

agtcta agcccg ctaatga gcgggct ttttttt ! terminator

GGCCcgactGGCC -3' ! Sfi I  
Sfi I



Table 103: Catalogue of plasmids

|    |         |   |
|----|---------|---|
| 5  | pEP1001 | pAA3H with 4.3 kbp deletion of $\lambda$ , <u>Cla</u> I site introduced.  |
| 10 | pEP1002 | pEP1001 with <u>fd</u> terminator and <u>Spe</u> I, <u>Sfi</u> I, <u>Hpa</u> I cloning site distal to <u>galT.K</u> .   |
| 15 | pEP1003 | pEP1002 with $P_1$ promoter replaced by pBR322 <u>amp</u> promoter ( <u>Pamp</u> ) and $O_R3$ upstream of <u>galT.K</u> ; <u>Pamp</u> and $O_R3$ bounded by <u>Hpa</u> I and <u>Xba</u> I and containing <u>Apa</u> I cloning site between <u>Hpa</u> I and <u>Pamp</u> . |
| 20 | pEP1004 | pKK175-6 with ( <u>Pamp</u> , <u>galT.K</u> , <u>fd</u> terminator, <u>Spe</u> I, <u>Sfi</u> I cloning site) from pEP1003   |
| 25 | pEP1005 | pEP1004 with Tn5 <u>neo</u> promoter ( <u>Pneo</u> ) and $O_R3$ bounded by <u>Stu</u> I and <u>Hind</u> III.  |
| 30 | pEP1006 | pEP1005 with <u>BamH</u> I site removed by site-specific mutation.  |
| 35 | pEP1007 | pEP1006 with ( <u>lacUV5</u> , S.D., <u>ray</u> cloning site, <u>trpa</u> terminator).  |
| 40 | pEP1008 | pEP1007 with N-terminal part of <u>ray</u> gene.  |
| 45 | pEP1009 | pEP1008 with complete <u>ray</u> gene.  |
|    | pEP1010 | pEP1009 with $O_R3$ replaced by scrambled $O_R3$ sequence.  |
| 50 | pEP1011 | pEP1009 with $O_R3$ sequences replaced with the HIV 353-369 Left Symmetrized Target.  |
| 55 | pEP1012 | pEP1009 With $O_R3$ sequences replaced with the HIV 353-369 Right Symmetrized Target.   |

Table 103 (continued) : Catalogue of plasmids

|    |                       |   |
|----|-----------------------|---|
| 5  | pEP1100 to<br>pEP1199 | pEP1011 with <u>ray<sub>L</sub></u> .   |
|    | pEP1200 to<br>pEP1299 | pEP1012 with <u>ray<sub>R</sub></u> .   |
| 10 | pEP1301               | pEP1100 with <u>ray<sub>L</sub></u> <sup>-</sup> VF55.  |
|    | pEP1302               | pEP1100 with <u>ray<sub>L</sub></u> <sup>-</sup> FW58.  |
| 15 | pEP1303               | pSP64 with Tn5 <u>neo</u>   |
|    | pEP1304               | pEP1303 with deletion of Ap resistance<br>gene.   |
| 20 |                       |   |
|    | pEP1305               | pEP1304 with <u>ray<sub>L</sub></u> <sup>-</sup> VF55.  |
| 25 | pEP1306               | pEP1304 with <u>ray<sub>L</sub></u> <sup>-</sup> FW58.  |
|    | pEP1307               | pEP1304 with <u>ray</u> .   |
| 30 | pEP1400 to<br>pEP1499 | pEP1200 series plasmids with HIV 353-369<br>substituted for Right Symmetrized Targets.  |
| 35 | pEP1500 to<br>pEP1599 | pEP1400 series plasmids containing<br>modified <u>ray<sub>R</sub></u> genes producing Ray <sub>R</sub> proteins<br>that complement the <u>ray<sub>L</sub></u> <sup>-</sup> VF55 mutation. |
| 40 | pEP1600 to<br>pEP1699 | pEP1400 series plasmids containing<br>modified <u>ray<sub>R</sub></u> genes producing Ray <sub>R</sub> proteins<br>that complement the <u>ray<sub>L</sub></u> <sup>-</sup> FW58 mutation. |
| 45 | pEP2000               | pEP1009 with <u>ray</u> replaced by <u>arc</u> .  |
|    | pEP2001               | pEP2000 with <u>arc</u> operator in <u>Pneo</u> , <u>tet</u> .  |
| 50 | pEP2002               | pEP2001 with <u>arc</u> operator in <u>Pamp</u> , <u>galT.K</u> .   |
| 55 |                       |   |

Table 103 (continued) : Catalogue of plasmids

|    |             |   |
|----|-------------|---|
| 5  | pEP2003     | pEP2002 with Target#1 in <u>Pneo</u> , <u>tet</u> .         |
|    | pEP2004     | pEP2003 with Target#1 in <u>Pamp</u> , <u>galT.K</u> .      |
| 10 | vg1-pEP2005 | pEP2004 with vgDNA (variegation #1 of polypeptide).         |
| 15 | vg2-pEP2006 | pEP2004 with vgDNA (variegation #2 of polypeptide).         |
| 20 | vg3-pEP2007 | pEP2004 with vgDNA (variegation #3 of polypeptide).         |
| 25 | pEP2010     | pEP2002 with Target#2 in <u>Pneo</u> , <u>tet</u> .         |
|    | pEP2011     | pEP2010 with Target#2 in <u>Pamp</u> , <u>galT.K</u> .      |
| 30 | vg1-pEP2012 | pEP2011 with vgDNA (variegation #1 of residues 1-10).       |
| 35 | pEP3000     | pEP2004 with <u>CI2-arc(1-10)</u> in place of <u>arc</u> .  |
|    | pEP4000     | pEP2002 with Target#3 in <u>Pneo</u> , <u>tet</u> .         |
| 40 | pEP4001     | pEP4000 with Target#3 in <u>Pamp</u> , <u>galT.K</u> .      |
|    | pEP4002     | pEP4001 with <u>cro-h12</u> in place of <u>arc</u> .        |
| 45 | vg1-pEP1233 | pEP4002 with vgDNA (variegation #1 of polypeptide segment). |
| 50 |             |   |
| 55 |             |   |

Table 104: fd t rminator  
and multiple cl ning site  
to insert after galt.X

5' |CGA|AAG|GCT|CCT|TTT|GCA|GCC|TTT|TTT|TTT|-  
Olig#2 = 3' t ttc cga gga aaa cgt cgg aaa aaa aaa|-  
|fd terminator|

|ACT|AGT|CAG|TGG|CCC|GAC|TGG|CCG|TTA|AC 3' = Olig#1  
|tga tca|gtc acc ggg ctg acc ggc aat tgg c 5'  
|SpeI| |SfiI| |HpaI|

Table 105: Mutag nic Primer  
to Remove BamHI site from pEP1005

5

10

```

      | t | p | v | l | w | i |
      | 93| 94| 95| 96| 97| 98|
5' CC|ACA|CCC|GTC|CTG|TGG|ATC|-
3' gg tgt ggg cag gac acc tat-

```

15

20

```

      | l | y | a | g | r | i |
      | 99|100|101|102|103|104|
      |CTG|TAC|GCC|GGA|CGC|ATC|GT 3' pEP1005
      Aac atg cgg cct gcg tag ca 5' Olig#7

```

25

**Bold, upper case bases indicate sites of mutation.**

30

35

40

**Table 106:     deleted.**

45

50

55

Table 107: Synthesis of lacUV5-BstEII-BglII-KpnI-trpA terminator

5' CTA|GTC|CAG|GCT|TTA|CAC|TTT|ATG|CTT|CCG|GCT|-  
 olig#9 = 3' ag gtc cga aat gtg aaa tac gaa ggc cga-  
SpeI -35

/3' = Olig#8  
CGT|ATA|ATG|TGT|GGA|ATT|GTG|AGC|GGA|TAA|CAA|TTT|-  
gtg tgt tac aca cct taa cac tcg cct att gtt aaa-  
-10 lac operator  
 Olig #11 = 3' /

/3' = Olig #10  
CAC|ACA|GGT|AAC|CAGGAGAGA TCT A|TGC|GGT|ACC|-  
gtg tgt cca ttg gtcctctct aga t acg cca tgg-  
BstEII BglII KpnI  
 Olig #13 = 3' /

AGT|CTA|AGC|CCG|CCT|AAT|GAG|CGG|GCT|TTT|TTT|TT-  
tca gat tcg ggc gga tta ctc gcc cga aaa aaa aa-  
spacer trpA terminator

G|GCC|CGA|C 3' = Olig #12  
c cga g 5'  
SfiI

Table 108: deleted.

Table 109: Synthesis of  
First segment of ray gene

5

10

5' C|AGG|AGG|TAA|CCA|gga|gga|aat|aaa|-  
|BstEII|

15

|ATG|GAA|CAA|CGC|ATA|ACC|CTA|AAG|GAC|TAC|GCG|ATG|CGC|-

20

/3' = Olig#14

|TTT|GGC|CAA|ACC|AAG|ACA|GCG|AAG|GAC|CTA|

Olig#15 = 3' gg ttc tgt cgc ttc ctg gat-

|BglI|

|PvuMI|

25

|GGG|GTG|TAT|CAG|AGC|GCG|ATT|AAC|AAG|GCC|ATC|

ccc cac ata gtc tcg cgc taa ttg ttc cgg tag-

30

|CAT|GCC|GGC|CGA|AAG|ATC|TTC|CTG|

gta cgg ccg gct ttc tag aag gac 5'

35

|BglII|

40

45

50

55

Table 110: Second segment of ray gene

5

| r | k | i | f | l | t | i | n | a | d |

| 38 | 39 | 40 | 41 | 42 | 43 | 44 | 45 | 46 | 47 |

10

5' C | CGA | AAG | ATC | TTC | CTA | ACC | ATT | AAC | GCT | GAT |

| BglII |

15

| g | s | v | y | a | e | e | v | k | p | f | p | s |

| 48 | 49 | 50 | 51 | 52 | 53 | 54 | 55 | 56 | 57 | 58 | 59 | 60 |

| GGA | TCC | GTC | TAC | GCG | GAA | GAG | GTA | AAG | CCC | TTC | CCG | AGT |

20

| BamHI || strand overlap || AvaI |

25

| n | k | k | t | t | a | . | . | . |

| 61 | 62 | 63 | 64 | 65 | 66 | 67 | 67 | 68 |

| AAC | AAA | AAA | ACA | ACA | GCG | TAA | TAG | TAG | gta | cca | gtc | t 3' |

30

| KpnI |

35

40

45

50

55



Table 111 :  $\lambda$  O<sub>R</sub> core sequences  
Used to search HIV-1

5

1234567

Kim et al. Consensus-A 5' CCGCGGG 3'

10

3' GGCGCCC 5' Kim Consensus-S

Symmetric Consensus-A 5' CCGCCGG 3'

15

3' GGCGGCC 5' Symm. Consensus-S

O<sub>R</sub>3A

5' CCGCAAG 3'

20

3' GGCGTTC 5' O<sub>R</sub>3SO<sub>R</sub>3A/Symm. Consensus.6 5' CCGCAGG 3'

25

3' GGCGTCC 5' O<sub>R</sub>3S/Symm. Cons.2O<sub>R</sub>3A/Symm. Consensus.5 5' CCGCCAG 3'

30

3' GGCGGTC 5' O<sub>R</sub>3S/Symm. Cons.3

7654321

35

40

45

50

55

Table 112: Potential target binding sequences  
 having subsequences matching  
 six of seven bases

5

10

|  
 CCGCGGG Kim consensus-A  
 HIV-1 subsequence =ACTTTCCGCTGGGGACT  
 353 ↑

15

20

|  
 CCGCAGG O<sub>R</sub>3A/consensus.6  
 HIV-1 subsequence =TCTCGACGCAGGACTCG  
 681 ↑

25

30

|  
 CTTGCGG O<sub>R</sub>3S  
 HIV-1 subsequence =TTTGACTAGCGGAGGCT  
 760 ↑

35

40

45

50

55

Table 113: Potential target binding sequences  
having subsequences matching five of seven bases

5

Symmetric consensus-S                    CCGGCGG  
HIV-1 subsequence    GACTTTCGGctGGGGAC  
352 ↑

10

OR3S/consensus.2                    CCTGCGG  
HIV-1 subsequence    TTTCCgCTGgGGACTTT  
355 ↑

15

OR3S/symm consensus.3                    CTGGCGG  
HIV-1 subsequence    TAGCAgTGGCGgCCGAA  
630 ↑

20

Symmetric consensus-A                    CCGCCGG  
HIV-1 subsequence    CAGTGgCGCCgGAACAG  
633 ↑

25

Or3A/symm consensus.5                    CCGCCAG  
HIV-1 subsequence    CAGTGgCGCCgGAACAG  
633 ↑

30

OR3A/consensus.6                    CCGCAGG  
HIV-1 subsequence    GACTAgCGgAGGCTAGA  
763 ↑

35

symm consensus-S                    CCGGCGG  
HIV-1 subsequence    GACTAgCGgAGGCTAGA  
763 ↑

40

45

50

55

Table 113, continued: Potential target binding sequences  
having subsequences matching five of seven bases

5

Or3A/sym consensus.5            CCGCCAG  
10       HIV-1 subsequence GAAGATgGCCAGTAAAA  
             4545 ↑

15

OR3A/consensus.6            CCGCAGG  
     HIV-1 subsequence ACAGATgGCAGGTGATG  
             5047 ↑

20

OR3A/consensus.6            CCGCAGG  
     HIV-1 subsequence TCCTATgGCAGGAAGAA  
             5965 ↑

25

30

35

40

45

50

55

**Table 114: Coding region of ray<sub>L</sub>-27 gene**

```

5      m   e   q   r   i   t   l   k   d   y   a   m   r
      1   2   3   4   5   6   7   8   9   10  11  12  13
10     ATG GAA CAA CGC ATA ACC CTA AAG GAC TAC GCG ATG CGC

      f   g   R   t   k   t   a   k   d   l
      14  15  16  17  18  19  20  21  22  23
15     TTT GGC CGT ACC AAG ACA GCG AAG GAC CTA
                                   |PpuM I|

      g   v   H   I   T   a   i   Q   N   a   i
      24  25  26  27  28  29  30  31  32  33  34
20     GGG GTG CAT ATT ACG GCG ATT CAG AAT GCC ATC

      h   a   g   K   Q   i   f   l   t   i   n   a   d
      35  36  37  38  39  40  41  42  43  44  45  46  47
25     CAT GCC GGC AAG CAG ATC TTC CTA ACC ATT AAC GCT GAT

      g   s   v   y   a   e   e   v   k   p   f   p   s
      48  49  50  51  52  53  54  55  56  57  58  59  60
30     GGA TCC GTC TAC GCG GAA GAG GTA AAG CCC TTC CCG AGT
      |BamHI|                               |Ava I|

      n   k   k   t   t   a   .   .   .
      61  62  63  64  65  66  67  67  68
40     AAC AAA AAA ACA ACA GCG TAA TAG TA GGTACC
                                   |KpnI|

```

Table 115: rayR-38 gene

5           m   e   q   r   i   t   l   k   d   y   a   m   r  
           1   2   3   4   5   6   7   8   9   10 11 12 13  
 ATG GAA CAA CGC ATA ACC CTA AAG GAC TAC GCG ATG CGC  
 10  
           f   g   E   t   k   t   a   k   d   l  
           14 15 16 17 18 19 20 21 22 23  
 TTT GGC GAG ACC AAG ACA GCG AAG GAC CTA  
 15                                   PpuM I  
  
           g   v   R   T   L   a   i   R   D   a   i  
 20           24 25 26 27 28 29 30 31 32 33 34  
 GGG GTG CGT ACT CTT GCG ATT CGT GAT GCC ATC  
  
           K   a   g   N   H   i   f   l   t   i   n   a   d  
 25           35 36 37 38 39 40 41 42 43 44 45 46 47  
 AAG GCC GGC AAT CAT ATC TTC CTA ACC ATT AAC GCT GAT  
  
 30  
           g   s   v   y   a   e   e   v   k   p   f   p   s  
           48 49 50 51 52 53 54 55 56 57 58 59 60  
 GGA TCC GTC TAC GCG GAA GAG GTA AAG CCC TTC CCG AGT  
 35                   BamHI                                   Ava I  
  
  
 40           n   k   k   t   t   a   .   .   .  
           61 62 63 64 65 66 67 67 68  
 AAC AAA AAA ACA ACA GCG TAA TAG TA GGTACC  
 45                                   KpnI  
  
  
 50  
  
  
 55

## Tables for Example 2

[0406]

Table 200

| P22 <u>arc</u> operator                         |                               |
|---|-------------------------------|
| P22 <u>arc</u> Operator                         | 5' ATGATAGAAG C ACTCTACTAT 3' |
|   | 3' TACTATCTTC G TGAGATGATA 5' |
| consensus of half-sites                         | 5' ATrTAGArk s myTCTAyyAT 3'  |
|   | 3' TAYyATCTym s krAGATrrTA 5' |
| P22 <u>arc</u> left half operator = ATrTAGArk   |                               |
| P22 <u>arc</u> right half operator = myTCTAyyAT |                               |

Table 201 P22 Arc gene

| m | k | g | m | s | k |  
 | 1 | 2 | 3 | 4 | 5 | 6 |  
 GG|TAA|CCT|ATG|AAG|GGT|ATG|TCT|AAA|-  
|BstE II|

| m | p | h | f | n | l | r | w | p | r |  
 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 |  
 |ATG|CCT|CAC|TTT|AAC|CTC|AGG|TGG|CCC|CGG|G-  
|Bsu36I|      |Xma I|

| e | v | l | d | l | v | r | k | v | a |  
 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 |  
 |AG|GTC|CTT|GAT|CTT|GTT|CGC|AAG|GTT|GCT|-  
|PpuM I|

| e | e | n | g | r | s | v | n | s | e |  
 | 27 | 28 | 29 | 30 | 31 | 32 | 33 | 34 | 35 | 36 |  
 |GAG|GAA|AAC|GGT|CGG|TCC|GTT|AAC|TCT|G|-  
|Rsr II|

|Hpa I|

Table 201, continued

5       | i | y | n | r | v | m | e | s | f | k |  
       | 37| 38| 39| 40| 41| 42| 43| 44| 45| 46|  
 AG|ATT|TAT|AAT|CGC|GTT|ATG|GAG|TCG|TTC|AAG|-  
 10   | Bgl II |

      | k | e | g | r | i | g | a | . | . | . |  
 15   | 47| 48| 49| 50| 51| 52| 53|   |   |   |  
       | AAA|GAG|GGT|CGT|ATC|GGC|GCA|TAA|TAG|TGA|

20       | GGT|ACC|  
       | Kpn I |

25       Amino acid sequence encoded is identical to wild type P22  
       Arc.

30       DNA sequence designed for optimal placement of restriction  
       sites.

35

40

45

50

55



10

3'- ga tac ttc cca tac aga ttt-

15

1ATG2CCT3CAC4TTT5AAC6CTC7AGG8TGG9CCC10CGG11-

20

olig#405

| GAG | GTC | CTT | GAT | CTT | GTT | CGC | AAG | GTT | GCT | -

ctc cag gaa cta gaa caa gcg ttc caa cga-

25

GAG GAA AAC GGT CGG TCC G TT AAC TCT GAG -

ctc ctt    ttg cca gcc agg c    aa ttg aga ctc-

\=olig#406

30

ATC TAT AAT CGC GTT ATG GAG TCG TTC AAG -

tag ata tta gcg caa tac ctc agc aag ttc-

\=olig#407

35

AAA GAG GGT CGT ATC GGC GCA TAA TAG TGA -

ttt ctc cca qca tag ccg cgt att atc act-

40

GGTAC 3' = olig#403

£ 5' = olig#408

1 Kpn I 1

49

Number of bases in each oligonucleotide.

50

**400 = 43**

**401 = 48**

**402 = 42**

403 - 47

**405 = 50**

**406 = 49**

**407 = 38**

**408 = 34**

Tabl 203: HIV-1 Subsequences  
that are similar to one half of  
the Arc Operator

5

10

Number of  
mismatches

15

1234567890|0987654321  
arco =ATrrTAGArk  
HIV-1 subsequence =ATtATAtAATACAGTAGCAAC 2  
1019 ↑

20

1234567890|0987654321  
arco =ATrrTAGArk  
HIV-1 subsequence =ATAATAgAGTAGCAACCCTCT 1  
1024 ↑

25

30

1234567890|0987654321  
arco = myTCTAyyAT  
HIV-1 subsequence =ACAGTAGCAACCCTCTATTgT 1  
1040 ↑

35

40

1234567890|0987654321  
arco =ATrrTAGArk  
HIV-1 subsequence =ATGATAGgGGGAATTGGAGGT 1  
2387 ↑

45

50

1234567890|0987654321  
arco =ATrrTAGArk  
HIV-1 subsequence =tTGAcAGAGAAAAAATAAAA 2  
2624 ↑

55

Table 204 Synthesis of Potential DBP-1  
vg1 for pEP2004

5

M K G M S K  
1 2 3 4 5 6

10

5'-GCCGTACGG|TAA|CCT|ATG|AAG|GGT|ATG|TCT|AAA|-  
|BstE II|

15

2 2 2 2 2 2  
M P Q F |I/M|Q/R|D/V|R/I|W/G|D/G|  
7 8 9 10| 11| 12| 13| 14| 15| 16|  
|ATG|CCT|CAC|TTT|ATs|CrG|GwT|AKA|KGG|GrT|-

20

|-3' = olig#420

25

2 2 2 2 2      ↓2 2 2  
|Q/L|R/T|F/Y|R/C|W/G| V | Q |I/M|T/I|R/Q|  
| 17| 18| 19| 20| 21| 22| 23| 24| 25| 26|  
|CWG|AsA|TWT|vGT|KGG|GTG|CAG|ATs|Ayc|CrG|  
3' -cc cac gtc taS tRg qYc-

30

35

2 2 2 2 2 2 2 2 2 2  
|V/I|R/I|F/Y|D/V|T/I|R/Q|V/I|D/G|V/I|P/Q|  
| 27| 28| 29| 30| 31| 32| 33| 34| 35| 36|  
|rTT|AKA|TWT|GwT|Ayc|CrG|rTT|GrT|rTT|CmG|  
Yaa tMt aWa cWa tRg qYc Yaa cYa Yaa qKc-

40

45

. . .  
|TAA|TAG|TGA|AAC|CTC|AGG|CGTGATCC  
att atc act ttg gag tcc gcactagg -5'=olig#421  
| Bsu36I | spacer|

50

55

Table 204, continued: NOTES

s = equimolar C and G                      r = equimolar A and G  
w = equimolar A and T                      k = equimolar T and G  
y = equimolar T and C                      m = equimolar A and C  
n = equimolar A, C, G, and T

There are  $2^{24} = (\text{approx.}) 1.6 \times 10^7$  DNA and protein sequences.

Number of bases in each oligonucleotide.

420 = 86                      421 = 73

Table 205: Result f first variegation

5

|   |   |   |   |   |   |
|---|---|---|---|---|---|
| M | K | G | M | S | K |
| 1 | 2 | 3 | 4 | 5 | 6 |

10

|ATG|AAG|GGT|ATG|TCT|AAA|-

15

|   |   |   |    |    |    |    |    |    |    |
|---|---|---|----|----|----|----|----|----|----|
| M | P | Q | F  | M  | R  | D  | I  | W  | G  |
| 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 |

|ATG|CCT|CAC|TTT|ATG|CGG|GAT|ATA|TGG|GGT|-

20

25

|    |    |    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|----|----|
| Q  | T  | Y  | C  | G  | V  | Q  | M  | T  | R  |
| 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 |

|CAG|ACA|TAT|TGT|GGG|GTG|CAG|ATG|ACC|CGG|-

30

35

|    |    |    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|----|----|
| V  | I  | F  | D  | I  | R  | V  | G  | V  | P  |
| 27 | 28 | 29 | 30 | 31 | 32 | 33 | 34 | 35 | 36 |

|GTT|ATA|TTT|GAT|ATC|CGG|GTT|GGT|GTT|CCG|

40

45

50

55

Table 206 Synthesis of Potential DBP14  
vg2 for pEP2004

5

M K G M S K  
1 2 3 4 5 6

10

5'-GCCGTACGG|TAA|CCT|ATG|AAG|GGT|ATG|TCT|AAA|-  
|BstE II|

15

M P Q F M|T R|Q D|N I|T W|R G|C  
7 8 9 10 11 12 13 14 15 16  
|ATG|CCT|CAC|TTT|Ayg|CrA|rAT|AvT|yGG|kGT|-

20

| = 3' olig#424

↓ V|D

25

Q|H T|N Y C G N|I Q|R M|T T|N R|C  
17 18 19 20 21 22 23 24 25 26  
|CAW|AMC|TAC|TGC|GGG|rwT|CrG|Ayg|Ayc|yGT|-

30

3' -g atg acg ccc YWa gYc tRc tRg Rca -  
| overlap |

35

V|P R|Q F|S D|N I|T R|I V|G G|R V|D P|R  
27 28 29 30 31 32 33 34 35 36  
|kTT|CrG|TyT|rAT|Ayc|AKA|GkT|sGT|GwT|CsG|  
Maa gYc aRa Yta tRg tMt cMa Sca cWa gSc -

40

TAA|TAG|TGA|AAC|CTC|AGG|CGTGATCC  
att atc act ttg gag tcc gcactagg -5'-olig#423  
| Bsu36I | spacer |

45

50

55

Tabl 206, c ntinu d: NOTES

5 s = equimolar C and G r = equimolar A and G  
 w = equimolar A and T k = equimolar T and G  
 y = equimolar T and C m = equimolar A and C  
 10 n = equimolar A, C, G, and T

15  $2^{24}$  sequences =  $1.6 \times 10^7$  sequences (DNA and protein).

Number of bases in each oligonucleotide.

20  $424 = 78$   $423 = 81$

25

30

35

40

45

50

55

Table 207 Result of second selection

5

M K G M S K

1 2 3 4 5 6

10

|ATG|AAG|GGT|ATG|TCT|AAA|-

15

M P Q F M R N I W G

7 8 9 10 11 12 13 14 15 16

|ATG|CCT|CAC|TTT|ATG|CGA|AAT|ATT|TGG|GGT|-

20

Q T Y C G D R M T R

17 18 19 20 21 22 23 24 25 26

25

|CAT|ACC|TAC|TGC|GGG|GAT|CGG|ATG|ACC|CGT|-

30

F N S N I R G R V R

27 28 29 30 31 32 33 34 35 36

|TTT|AAT|TCT|AAT|ATC|AGA|GGT|CGT|GTT|CGG|

35

40

. . .  
|TAA|TAG|TGA|

45

50

55



Tabl 208: Third vari gation vg3 for pEP2004

5 M K G M S K  
1 2 3 4 5 6

5'- CGTCGCATGG|TAA|CCT|ATG|AAG|GGT|ATG|TCT|AAA|-  
|spacer|BstE II|

10 M|K N|D W|S  
M P Q F E|V R T|A I R|P G  
7 8 9 10 11 12 13 14 15 16

15 |ATG|CCT|CAC|TTT|rwG|CGG|rwT|ATA|vsG|GGT|-

20 Q|R Y|C D|I R|Q  
G|E T H|R C G N|V G|E M T R  
17 18 19 20 21 22 23 24 25 26

|srG|ACA|vrT|TGT|GGG|rwT|srG|ATG|ACC|CGC|-olig#325  
olig#327 3'- c tac tgg gcg-  
25 | overlap |

30 F|C S|N I|T G|D R|H  
V|G N R|H N P|L R R|H R V P|L  
27 28 29 30 31 32 33 34 35 36

|kkT|AAT|mrT|AAT|myC|CGG|srT|CGT|GTT|Cnt|  
35 Mma tta KYa cta KRg gtc SYa gca caa gNa-

40 TAA|TAG|TGA|AAC|CTC|AGG|CGACCTGGC  
att atc act ttg gag tcc gctggaccg -5'  
| Bsu36I |

45 s = equimolar C and G r = equimolar A and G  
w = equimolar A and T k = equimolar T and G  
y = equimolar T and C m = equimolar A and C  
50 n = equimolar A, C, G, and T

$$4^{12} = 2^{24} = 1.6 \times 10^7 \text{ protein and DNA sequences}$$

Table 209: Polypeptide that  
Binds First Targ t

5

10

|                           |   |   |   |   |   |
|---------------------------|---|---|---|---|---|
| M                         | K | G | H | S | K |
| 1                         | 2 | 3 | 4 | 5 | 6 |
| ATG AAG GGT ATG TCT AAA - |   |   |   |   |   |

15

20

|   |   |   |    |    |    |    |    |    |    |
|---|---|---|----|----|----|----|----|----|----|
| M   | P | Q | F  | V  | R  | D  | I  | R  | G  |
| 7   | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 |
| ATG CCT CAC TTT GTG CGG GAT ATA CGG GGT - |   |   |    |    |    |    |    |    |    |

25

30

|   |    |    |    |    |    |    |    |    |    |
|---|----|----|----|----|----|----|----|----|----|
| G                                       | T  | H  | C  | G  | I  | Q  | M  | T  | R  |
| 17                                      | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 |
| GGG ACA CAT TGT GGG ATT CAG ATG ACC CGC |    |    |    |    |    |    |    |    |    |

35

40

|   |    |    |    |    |    |    |    |    |    |
|---|----|----|----|----|----|----|----|----|----|
| V                                       | N  | R  | N  | P  | R  | H  | R  | V  | L  |
| 27                                      | 28 | 29 | 30 | 31 | 32 | 33 | 34 | 35 | 36 |
| ATT AAT CGT AAT CCC CGG CAT CGT GTT CTT |    |    |    |    |    |    |    |    |    |

45

50

55

Table 210: Variegation for Second Target vgl for pEP2011

5 K|T  
A|V  
M|T G|D G|D M|T S|R K|Q  
10 M V|A R|M A|V V|A N|K N|H  
0 1 2 3 4 5 6  
5'- CGTCGCATGG|TAA|CCT|ATG|ryG|rNG|GnT|ryG|Ars|mas|-  
15 |spacer|BstE II|  
F|Y  
H|L  
20 M|K P|Q Q|H I|N V|I  
Q|L R|L N|A V|D T|A R D I R G  
7 8 9 10 11 12 13 14 15 16  
25 |mwG|CnT|mas|nWT|ryT|CGG|GAT|ATA|CGG|GGT|-  
/ = olig#460  
30 G T H C G I | Q M T R  
17 18 19 20 21 22 | 23 24 25 26  
|GGG|ACA|CAC|TGC|GGG|ATC| CAG|ATG|ACC|CGC|  
olig#461 = 3'-gtg acg ccc tag gtc tac tgg acg-  
35 | overlap |  
V N R N P R H R V L  
40 27 28 29 30 31 32 33 34 35 36  
|ATT|AAT|CGT|AAT|CCC|CGG|CAT|CGT|GTT|CTT|  
taa tta gca tta ggg gcc gta gca caa gaa  
45  
TAA|TAG|TGA|AAC|CTC|AGG|CGACCTGGC -3'  
50 att atc act ttg gag tcc gctggaccg -5'  
| Bsu36I | spacer |  
55

Table 210, continued: NOTES

s = equimolar C and G                      r = equimolar A and G  
 w = equimolar A and T                      k = equimolar T and G  
 y = equimolar T and C                      m = equimolar A and C  
 n = equimolar A, C, G, and T

$2^{24} = 1.6 \times 10^7$  protein and DNA sequences

Table 211: Polypeptide Selected  
for Binding to Second Target

|                               |   |   |   |   |   |   |
|-------------------------------|---|---|---|---|---|---|
| M                             | T | R | D | M | K | Q |
| 0                             | 1 | 2 | 3 | 4 | 5 | 6 |
| ATG ACG AGG GAT ATG AAG CAG - |   |   |   |   |   |   |

|   |   |   |    |    |    |    |    |    |    |
|---|---|---|----|----|----|----|----|----|----|
| M   | Q | N | D  | I  | R  | D  | I  | R  | G  |
| 7   | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 |
| ATG CAT AAC GAT ATT CGG GAT ATA CGG GGT - |   |   |    |    |    |    |    |    |    |

|   |    |    |    |    |    |    |    |    |    |
|---|----|----|----|----|----|----|----|----|----|
| G                                       | T  | H  | C  | G  | I  | Q  | M  | T  | R  |
| 17                                      | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 |
| GGG ACA CAC TGC GGG ATC CAG ATG ACC CGC |    |    |    |    |    |    |    |    |    |

|   |    |    |    |    |    |    |    |    |    |
|---|----|----|----|----|----|----|----|----|----|
| V                                       | N  | R  | N  | P  | R  | H  | R  | V  | L  |
| 27                                      | 28 | 29 | 30 | 31 | 32 | 33 | 34 | 35 | 36 |
| ATT AAT CGT AAT CCC CGG CAT CGT GTT CTT |    |    |    |    |    |    |    |    |    |

Tables for Example 3  
Table 300: CI2-arc(1-10) gene

5

10

15

20

25

30

35

40

45

50

55

|  |  |
|--|--|
|  | m   l   k   t   e   w                                |
|  | 1   2   3   4   5   6                                |
|  | GG TAA CCT ATG CTT AAG ACT GAA TGG                   |
|  | <u>BstEII</u>     <u>Afl I</u>                       |
|  | p   e   l   v   g   k   s   v   e   e                |
|  | 7   8   9   10   11   12   13   14   15   16         |
|  | CCT GAG CTT GTT GGT AAA TCT GTC GAG GAA              |
|  | a   k   k   v   i   l   q   d   k   p   e            |
|  | 17   18   19   20   21   22   23   24   25   26   27 |
|  | GCT AAG AAA GTT ATC CTG CAG GAT AAA CCT GAG          |
|  | <u>Pst I</u>     <u>Bsu36I</u>                       |
|  | a   q   i   i   v   l   p   v   g                    |
|  | 28   29   30   31   32   33   34   35   36           |
|  | GCC CAA ATC ATA GTA CTT CCG GTT GGC                  |
|  | <u>Sca I</u>   |
|  | t   i   v   t   m   e   y   r   i   d                |
|  | 37   38   39   40   41   42   43   44   45   46      |
|  | ACT ATT GTT ACC ATG GAG TAT CGT ATT GAC              |
|  | <u>Nco I</u>   |
|  | <u>Sty I</u>   |
|  | r   v   r   l   f   v   d   k   l   d                |
|  | 47   48   49   50   51   52   53   54   55   56      |
|  | CGC GTT CGT CTT TTT GTC GAC AAA TTG GAT              |
|  | <u>Acc I</u>   |
|  | <u>Hind II</u>                                       |
|  | <u>Sal I</u>   |

(continued on next page)

Table 300, continued

5 | n | i | a | e | v | p | r | v | g | g |  
 | 57 | 58 | 59 | 60 | 61 | 62 | 63 | 64 | 65 | 66 |  
 | AAC | ATT | GCT | GAG | GTC | CCT | CGC | GTA | GGT | GGC |

10 | Dra II |  
 | PvuM I |  
 | Pss I |  
 | Ava II |

15 | k | m | k | g | m | s | k | m | p | q |  
 | 67 | 68 | 69 | 70 | 71 | 72 | 73 | 74 | 75 | 76 |  
 20 | AAA | ATG | AAA | GGT | ATG | TCT | AAG | ATG | CCG | CAA |

| f | . | . | . |  
 25 | 77 | 78 | 79 | 80 |  
 | TTT | TAA | TGA | TAG | GGT | ACC |  
 | Asp718 |  
 | Kpn I |

30

Residue M1 is inserted so that translation can initiate.

35

Residue L2 corresponds to residue L20 of Barley chymotrypsin inhibitor CI-2.

40

Residues G66 and K67 are inserted to allow flexibility between CI-2 and the DNA-binding tail.

45

Residues 68-77 have the same sequence as the first ten residues of P22 Arc.

50

55

Table 301: Synthesis of  
CI2-arc(1-10) gene

5

10

15

20

25

30

| m | l | k | t | e | w |  
 | 1 | 2 | 3 | 4 | 5 | 6 |  
 5'-G|TAA|CCT|ATG|CTT|AAG|ACT|GAA|TGG|-  
 3'-ga tac gaa ttc tga ctt acc -  
 | BstEII | | Afl I |

/3' = olig#470  
 | p | e | l | v | g | k | s | | v | e | e |  
 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | | 14 | 15 | 16 |  
 CCT|GAG|CTT|GTT|GGT|AAA|TCT| GTC|GAG|GAA|-  
 gga ctc gaa caa cca ttt aga caa ctc ctt -

| a | | k | k | v | i | l | q | d | k | p | e |  
 | 17 | | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 |  
 GCT| AAG|AAA|GTT|ATC|CTG|CAG|GAT|AAA|CCT|GAG|-  
 cga ttc ttt caa tag gac gtc cta ttt gga ctc -  
 5' ↑ | Pst I | | Bsu36I |

olig#475

35

40

/3' = olig#471  
 | a | q | i | i | v | | l | p | v | g | |  
 | 28 | 29 | 30 | 31 | 32 | | 33 | 34 | 35 | 36 | |  
 GCC|CAA|ATC|ATA|GTA| CTT|CCG|GTT|GG C|-  
 cgg gtt tag tat cat gaa ggc caa cc g -  
 | Sca I | | 5' Olig#476

45

50

| t | i | v | t | m | e | y | r | i | d |  
 | 37 | 38 | 39 | 40 | 41 | 42 | 43 | 44 | 45 | 46 |  
 ACT|ATT|GTT|ACC|ATG|GAG|TAT|CGT|ATT|GAC|-  
 tga taa caa tgg tac ctc ata gca taa ctg  
 | Nco I | 5' olig#477 ↑  
 | Sty I |

55

Table 301, continued

5

/ 3' olig#472

|   |   |   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|---|---|
| r | v | r | l | f | v | d | k | l | d |
|---|---|---|---|---|---|---|---|---|---|

|    |    |    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|----|----|
| 47 | 48 | 49 | 50 | 51 | 52 | 53 | 54 | 55 | 56 |
|----|----|----|----|----|----|----|----|----|----|

10

|     |     |     |     |     |     |     |     |     |     |   |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|---|
| CGC | GTT | CGT | CTT | TTT | GTC | GAC | AAA | TTG | GAT | - |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|---|

|     |     |     |     |     |     |     |     |     |     |   |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|---|
| gcg | caa | gca | gaa | aaa | cag | ctg | ttt | aac | cta | - |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|---|

|       |
|-------|
| Acc I |
|-------|

|         |
|---------|
| Hind II |
|---------|

15

|       |
|-------|
| Sal I |
|-------|

olig#473 3' ↓

20

|   |   |   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|---|---|
| n | i | a | e | v | p | r | v | g | g |
|---|---|---|---|---|---|---|---|---|---|

|    |    |    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|----|----|
| 57 | 58 | 59 | 60 | 61 | 62 | 63 | 64 | 65 | 66 |
|----|----|----|----|----|----|----|----|----|----|

|     |     |     |     |     |     |     |     |     |     |   |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|---|
| AAC | ATT | GCT | GAG | GTC | CCT | CGC | GTA | GGT | GGC | - |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|---|

25

|     |     |     |     |     |     |     |     |     |     |   |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|---|
| ttg | taa | cga | ctc | cag | gga | gcg | cat | cca | ccg | - |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|---|

|        |   |
|--------|---|
| Dra II | ↑ |
|--------|---|

↑ 3' olig#479

|        |  |
|--------|--|
| PvuM I |  |
|--------|--|

|       |               |
|-------|---------------|
| Pss I | = 5' olig#478 |
|-------|---------------|

30

|        |
|--------|
| Ava II |
|--------|

35

|   |   |   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|---|---|
| k | m | k | g | m | s | k | m | p | q |
|---|---|---|---|---|---|---|---|---|---|

|    |    |    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|----|----|
| 67 | 68 | 69 | 70 | 71 | 72 | 73 | 74 | 75 | 76 |
|----|----|----|----|----|----|----|----|----|----|

|     |     |     |     |     |     |     |     |     |     |   |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|---|
| AAA | ATG | AAA | GGT | ATG | TCT | AAG | ATG | CCG | CAA | - |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|---|

40

|     |     |     |     |     |     |     |     |     |     |   |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|---|
| ttt | tac | ttt | cca | tac | aga | ttc | tac | ggc | att | - |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|---|

45

|   |   |   |   |
|---|---|---|---|
| f | . | . | . |
|---|---|---|---|

|    |    |    |    |
|----|----|----|----|
| 77 | 78 | 79 | 80 |
|----|----|----|----|

|     |     |     |     |     |      |               |
|-----|-----|-----|-----|-----|------|---------------|
| TTT | TAA | TGA | TAG | GGT | AC - | 3' = olig#474 |
|-----|-----|-----|-----|-----|------|---------------|

|     |     |     |     |   |      |
|-----|-----|-----|-----|---|------|
| aaa | att | act | atc | c | - 5' |
|-----|-----|-----|-----|---|------|

50

|       |
|-------|
| Kpn I |
|-------|



Table 301, continu d

Number of bases in each oligonucleotide.

5

|          |         |          |         |
|----------|---------|----------|---------|
| olig#470 | .... 46 | olig#475 | .... 53 |
| olig#471 | .... 57 | olig#476 | .... 56 |
| olig#472 | .... 54 | olig#477 | .... 31 |
| olig#473 | .... 48 | olig#478 | .... 48 |
| olig#474 | .... 47 | olig#479 | .... 55 |

10

15

20

25

30

35

40

45

50

55

Table 302: Variegation of Tail on CI-2 vgl for pEP3000

5' | e | v | p | r | v | g | g |  
 | 60 | 61 | 62 | 63 | 64 | 65 | 66 |  
cga gtc ggc | GAG | GTC | CCT | CGC | GTA | GGT | GGC |  
 | spacer | PpuM I |  
 3'

| k | m | k | g | m | s | k | m | p | q |  
 | 67 | 68 | 69 | 70 | 71 | 72 | 73 | 74 | 75 | 76 |  
AAA | ATG | AAA | GGT | ATG | AGC | AAG | ATG | CCG | CAG | -

| f | I/M | Q/R | D/V | R/G | G | V |  
 77 | 78 | 79 | 80 | 81 | 82 | 83 |  
TTC | ATs | CrG | GwT | sGA | GGT | GTC | -  
 olig#481 3'- ct cca caa -

/ 3' = olig#480  
 | Q/L | R/T | F/Y | R/C | W/G | V/D | Q/R | I/M | T/I | R/Q |  
 | 84 | 85 | 86 | 87 | 88 | 89 | 90 | 91 | 92 | 93 |  
C | Wg | AsA | TwT | yGT | kGG | GwC | CrG | ATs | AyC | CrG | -  
g | Wc | tSt | aWa | Rca | mCc | cWc | gYc | taS | tRg | gYc -

| V/I | R/I | F/Y | D/V | T/I | R/Q | V/I | D/G | V/I | P/Q |  
 | 94 | 95 | 96 | 97 | 98 | 99 | 100 | 101 | 102 | 103 |  
rTT | AKA | TwT | GwT | AyC | CrG | rTT | GrT | rTT | CmG |  
Yaa | tMt | aWa | cWa | tRg | gYc | Yaa | cYa | Yaa | gKc -

| . | . | . |  
 | 104 | 105 | 106 |  
 | TAA | TGA | TAG | GGT | AC  
att act atc c - 5'  
 | Kpn I |

2<sup>24</sup> DNA and protein sequences = (approx)  $1.6 \times 10^7$ .

Number of bases in each oligonucleotide.

480 ... 82

481 ... 78

## Tables f r Example 4

Table 400 Search for  $\lambda$  Cro Half Site in  
HIV non-variabl regions.

Gene seq id=HIVHXB2CG

## Sequences sought

|         | Consensus     | Or3-          | Or3-/consensus<br>hybrid |
|---------|---------------|---------------|--------------------------|
| forward | 5' TATCACC 3' | 5' TATCCCT 3' | 5' TATCACT 3'            |
| reverse | 5' GGTGATA 3' | 5' AGGGATA 3' | 5' AGTGATA 3'            |

## Match with Or3-

matches =TATCCCT  
HIV subsequence =aATCtCTAGCAGTGGCG  
624 ↑

## Match with Or3/consensus hybrid

matches =TATCACT  
HIV subsequence =aATCtCTAGCAGTGGCG  
624 ↑

## Match with Consensus

matches =GGTGATA  
HIV subsequence =ACAGATGGCAGGTGATg  
5057 ↑

Table 400, c ntinued

5

Match with Consensus

matches =TATCACC

10

HIV subsequence =CATCtCCTATGGCAGGA

5961 ↑

15

First target : TATCCCTAGCAGTGGCG

Second target: aATCtCTAGCAGTGGCG

20

624 ↑

25

30

35

40

45

50

55

Table 401: λ Cro alpha 1 & 2 & slot for Polypeptide

5' cga|CGG|AGG|TAA|CCT|ATG|GAA|CAA|CGC|ATA|ACC|-  
| spacer | BstE II |

olig#483 3'↓

15 | 1 | k | d | y | a | m | r | f | g | q |  
| 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 |  
| CTA | AAG | GAC | TAC | GCG | ATG | CGC | TTT | GGC | CAA |  
gc tac gcg aaa ccg gtt-

20 | Bal I |

25 | t | k | t | a | k | d | l | g | v |  
| 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 |  
| ACC | AAG | ACA | GCC | AAA | GAT | CTC | GGG | GTG |  
tgg ttc tgt cgg ttt cta gag ccc cac-

30 | Bgl II |  
| Ava I |

35 | . | . | . |  
| | | |  
| TAG | TAG | TAG | GGT | ACC | AAG | GCG |  
atc atc atc cca tgg ttc cgc - 5' olig#484  
40 | Kpn I | sapcer |

Number of bases in each oligonucleotide.

45

483 ... 60

484 ... 65

Table 402 Variegated Polypeptide to attach to  
Cro Helices 1, 2, & 3 vgl for pEP4002

5

k | d | l | g | v |

21| 22| 23| 24| 25|

ccg|acg|gcc|cga|GAT|CTC|GGG|GTG|-

10

| spacer | Bgl II |

| Ava I |

15

| y | q | s | a | i | n | k | a | i | h |

| 26| 27| 28| 29| 30| 31| 32| 33| 34| 35|

| TAT|CAG|AGC|GCG|ATT|AAC|AAA|GCG|ATC|CAC|-

20

I|M Q|R D|V R|I W|G D|G Q|L R|T F|Y R|C

| 36| 37| 38| 39| 40| 41| 42| 43| 44| 45|

| ATs|CrG|GWT|AKA|KGG|GrT|CwG|AsA|TWT|VGT|-

25

↓ 3' = olig#486

W|G V Q I|M T|I R|Q V|I R|I F|Y D|V

| 46| 47| 48| 49| 50| 51| 52| 53| 54| 55|

| KGG|GTG|CAG|ATs|Ayc|CrG|rTT|AKA|TWT|GWT|

30

cc cac gac taS yRg gYc Yaa tMt aWa cAa-

35

T|I R|Q V|I D|G V|I P|Q

| 56| 57| 58| 59| 60| 61|

| Ayc|CrG|rTT|GrT|rTT|CmG|

tRg gYc Yaa cYa Yaa gKc-

40

| . | . | . |

| | | |

| TAG|TAG|TAG|GGT|ACC|AAG|GCG|

atc atc atc cca tgg ttc cgc 5' = olig#488

45

| Kpn I | spacer |

50

s = equimolar C and G

r = equimolar A and G

w = equimolar A and T

k = equimolar T and G

y = equimolar T and C

m = equimolar A and C

n = equimolar A, C, G, and T

55

Table 403  
Result of first variation of  
alpha 1,2,3:vgPolyPeptide

5

10

15

20

25

30

35

40

45

50

55

|     |     |     |     |     |     |
|-----|-----|-----|-----|-----|-----|
| m   | e   | q   | r   | i   | t   |
| 1   | 2   | 3   | 4   | 5   | 6   |
| ATG | GAA | CAA | CGC | ATA | ACC |

|     |     |     |     |     |     |     |     |     |     |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| 1   | k   | d   | y   | a   | m   | r   | f   | g   | q   |
| 7   | 8   | 9   | 10  | 11  | 12  | 13  | 14  | 15  | 16  |
| CTA | AAG | GAC | TAC | GCG | ATG | CGC | TTT | GGC | CAA |

|     |     |     |     |     |     |     |     |     |     |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| t   | k   | t   | a   | k   | d   | l   | g   | v   | y   |
| 17  | 18  | 19  | 20  | 21  | 22  | 23  | 24  | 25  | 26  |
| ACC | AAG | ACA | GCC | AAG | GAC | CTA | GGC | GTG | TAT |

|     |     |     |     |     |     |     |     |     |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| q   | s   | a   | i   | n   | k   | a   | i   | h   |
| 27  | 28  | 29  | 30  | 31  | 32  | 33  | 34  | 35  |
| CAG | AGC | GCG | ATT | AAC | AAA | GCG | ATC | CAC |

|     |     |     |     |     |     |     |     |     |     |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| M   | Q   | V   | R   | G   | D   | L   | T   | Y   | C   |
| 36  | 37  | 38  | 39  | 40  | 41  | 42  | 43  | 44  | 45  |
| ATG | CAG | GTT | AGA | GGG | GAT | CTG | ACA | TAT | TGT |

|     |     |     |     |     |     |     |     |     |     |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| W   | V   | Q   | I   | I   | R   | V   | R   | F   | D   |
| 46  | 47  | 48  | 49  | 50  | 51  | 52  | 53  | 54  | 55  |
| TGG | GTG | CAG | ATC | ATC | CGG | GTT | AGA | TTT | GAT |

|     |     |     |     |     |     |
|-----|-----|-----|-----|-----|-----|
| T   | R   | V   | G   | I   | Q   |
| 56  | 57  | 58  | 59  | 60  | 61  |
| ACC | CGG | GTT | GGT | ATT | CAG |

|     |     |     |
|-----|-----|-----|
| .   | .   | .   |
|     |     |     |
| TAG | TAG | TAG |

## Tables for Example 5

Table 500: Proposed binding of Arc dimer to arcO.

(a) Interaction of residues 1-10  
with arcO

Arc            N----->C C<-----N  
arcO        5' ATrrTAGArksmyTCTAyyAT  
             3' TAYyATCTymskrAGATrrTA

(b) N-terminal residues interacting with same  
polypeptide chain, dimer contacts near C-terminus

                         2 C    C 1  
                         /    \  
                     /VVVV\VVVV\  
                     /VVVV\VVVV\  
Arc            1 N-----|-----N 2  
arcO        5' ATrrTAGArksmyTCTAyyAT  
             3' TAYyATCTymskrAGATrrTA

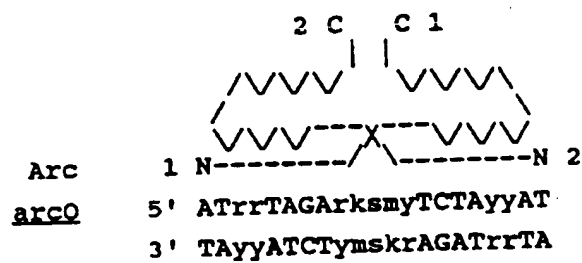


Table 500, continued

5 (c) N-terminal residues interacting with opposite  
 polypeptide chain, dimer contacts close to residue 10.

10

15



20

25

30

35

40

45

50

55

Table 501 : Search of HIV-1 isolate HXB2 DNA sequenc  
for sequences related to one half of arcO

In arcO sequence, upper case letters represent palindromical-  
ly related bases.

In HIV-1 subsequences, "@" represents a nucleotide found to  
vary among HIV-1 isolates while lower case letters represent  
mismatch to arcO.

HIV-1 1016-1051 is non-variable.

arcO left half = ATrrTAGArk  
HIV-1 subsequence =@@ATCATTATATAATACAGTAGCAACCCTCTATTGTGT@  
1024 ↑

arcO right half = myTCTAyyAT  
HIV-1 subsequence =CAGTAGCAACCCTCTATTgTGT@  
1040 ↑

2387-2427 is non-variable.

arcO left half = ATrrTAGArk  
HIV-1 subsequence =@ATGATAGgGGGAATTGGAGGTTTTATCAAAG  
2387 ↑

4661-4695 is non-variable.

arcO left half = ATrrTAGArk  
HIV-1 subsequence =AAGTCAAGGAgTAGTAGAATCTATGAATAA@  
4676 ↑

Table 502: Progression of Targets Leading to  
HIV-1 1016-1037

(a)

|                        |  |      |        |      |
|------------------------|--|------|--------|------|
|                        | 1016   | 1024 | center | 1047 |
|                        | ↓  | ↓    | ↓      | ↓    |
| HIV-1 5'               | ATCATTATATAATAcAGTAGCAACCCTCTATT                   |      |        |      |
| <u>First target</u> 5' | <u>TAcATgATAgAagcaCTataCTaTeee</u>                 |      |        |      |
| P22 sequence 5'        | <u>attgacATgaTAGAagcacTCTActATattctcaata</u> 3'    |      |        |      |
|                        | 3' <u>taactcTActATCTTcgtgAGATgaTataagagttat</u> 5' |      |        |      |
|                        | ----- <u>arcO</u> -----                            |      |        |      |

In target:      Upper case indicates that HIV-1 and arcO  
agree.  
Lower case indicates a change to match arcO.  
Underscore indicates identity to arcO.  
@ indicates bases that vary between instances  
of target.

In arcO :      underscore indicates DNase I protected.  
lower case indicates not palindromically  
related.

Table 502, continu d

(b)

Novel DBP = NXXXXX----->C C<-----XXXXXN

.....| | | | | | | | | | | | | # | | | | | \$\$\$

First target = @@@TACATgATAgAagcaCtAtaCTaT@@@

In the Novel DBP, X represents variegated sequence.

**In the line between Novel DBP and target DNA:**

represents regions where variegated sequence will produce amino acid sequences that will bind specifically.

N & C are the amino and carboxy ends residues 1-10.

| or \ represent regions where constant amino acid sequence is known to bind DNA.

# represents regions where constant amino acid sequence is believed not to bind DNA.

\$ represents regions where DNA sequence varies between different instances of the target.

Table 502: Progressi n of Targets Leading to  
HIV-1 1016-1034  
(continued)

(c)

|                         | 1016  | 1024  | center | 1047  |
|-------------------------|---|-------|--------|-------|
|                         | ↓   | ↓     | ⊕      | ↓     |
| HIV-1 5'                | ATCATTATATAATACAGTAGCAACCCTCTATT                |       |        |       |
| First target 5'         | <u>TACATgATAgAagcaCtAtaCTaT</u>                 |       |        |       |
| changes                 |   | ↓ ↓ ↓ |        | ↓ ↓ ↓ |
| <u>Second target</u> 5' | <u>TACATAATACAGgcaCtAtCCTee</u>                 |       |        |       |
| P22 sequence 5'         | <u>attgacATgaTAGAagcacTCTActATattctcaata</u> 3' |       |        |       |
| 3'                      | <u>taactgTActATCTtcgtgAGATgaTataagagttat</u> 5' |       |        |       |
|                         | ----- arco -----                                |       |        |       |

25

(d)

|                        |  |
|------------------------|--|
| Novel DBP =            | N---XXXXXXXXX>C C<XXXXXXXXX---N        |
|                        | \   .....   \$\$\$\$\$\$\$\$\$\$\$\$\$ |
| <u>Second target</u> = | @TACATAATACAGgcaCtAtCCTeeeeeee         |

35

40

45

50

55

Table 502, continued

5 (e) 1016 1024 center 1047  
 ↓ ↓ ↓ ↓  
 HIV-1 5' ATCATTATATAATACAGTAGCAACCCTCTATT  
 10 Second target 5' @@@TACATAATACAGGcaCtAtCCT  
 changes ↓↓↓ ↓ ↓↓↓ ↓ ↓↓↓  
 Third target 5' CATTATATAATACAGTAaCAACC@@  
 15 P22 sequence 5' attgacATgaTAGAagcactCTActATattctcaata 3'  
 3' taactgTActATCTtcgtgAGATgaTataagagttat 5'  
 |----- arco -----|

20

25

(f)

diffuse variegation

Novel DBP = NXXXXXXXXXXXXXXXXX&gt;C C&lt;XXXXXXXXXXXXXXXXN

.....| | \$\$\$\$\$\$\$\$\$\$\$\$\$\$\$\$\$\$

Third target = @@CATTATATAATACAGTAaCAACC@

30

35

40

45

50

55

10

15

20

25

30

35

10

15

22

Table 503: First Target Downstream  
of Promoters of Selectable Genes

First Target downstream  
of Pamp that promotes galT.K

5' |CCT|GCG|AAC|CGG|AAT|TGC|CAG|-  
Olig#501 = 3' gga cgc ttg gcc tta acg gtc-  
|StuI| | -35 |

|CTG|GGG|CGC|CCT|CTG|GTA|AGG|TTG|GGA|-  
gac ccc gcg gga gac cat tcc aac cct -  
| -10 |

1024

↓

|TAC|ATG|ATA|GAA|GCA|CTA|TAC|TAT|A 3' = olig#502  
atg tac tat ctt cgt gat atg ata t tca a 5'  
|First Target| |Hind3|



5

10

15

20

25

30

35

40

45

50

55

Table 504 First variegated  
insert into ped gen

5' | m | X | X | X | X | X |  
 | 1 | 96 | 97 | 98 | 99 | 100 |  
 | cct | cag | cGG | TAA | CCT | ATG | fzk | fzk | fzk | fzk | fzk |  
 | spacer | BstE II |  
 | m | k | g | m | s | k |  
 | 101 | 102 | 103 | 104 | 105 | 106 |  
 | ATG | AAG | GGT | ATG | TCT | AAA |  
 | m | p | h | f | n | l | r | Center of symmetry for  
 | 107 | 108 | 109 | 110 | 111 | 112 | 113 | ↓ priming.  
 | ATG | CCT | CAC | TTT | AAC | CTC | AGG | cgt | att | aat | acg | cct | g-3'  
 | Bsu36I | | primer |  
 olig#605 ↑

### Self priming

.....CTC|AGG|cgt|att|\  
 3'- tcc gca taa /

3' end self primes for extension with Klenow enzyme.

f = (0.26 T, 0.18 C, 0.26 A, 0.30 G)  
 z = (0.22 T, 0.16 C, 0.40 A, 0.22 G)  
 k = equimolar T and G

There are  $(2^5)^5 = 3.2 \times 10^7$  different DNA sequences encoding  
 $20^5 = 3.2 \times 10^6$  different prptein sequences.

100 has been added to residue numbers for wild-type Arc.

Table 505: Protein Ped-6  
Selected for Binding to First Target

5

| m | k | d | i | w | r |  
| 1 | 96 | 97 | 98 | 99 | 100 |  
| ATG | AAG | GAT | ATT | TGG | CGT |

10

| m | k | g | m | s | k |  
| 101 | 102 | 103 | 104 | 105 | 106 |  
| ATG | AAG | GGT | ATG | TCT | AAA |

15

| m | p | h | f | n | l | r | w | p | r |  
| 107 | 108 | 109 | 110 | 111 | 112 | 113 | 114 | 115 | 116 |  
| ATG | CCT | CAC | TTT | AAC | CTC | AGG | TGG | CCC | CGG | G |  
| Bsu36I | | Xma I |

20

| e | v | l | d | l | v | r | k | v | a |  
| 117 | 118 | 119 | 120 | 121 | 122 | 123 | 124 | 125 | 126 |  
| AG | GTC | CTT | GAT | CTT | GTT | CGC | AAG | GTT | GCT |  
| PvuM I |

25

| e | e | n | g | r | s | v | n | s | e |  
| 127 | 128 | 129 | 130 | 131 | 132 | 133 | 134 | 135 | 136 |  
| GAG | GAA | AAC | GGT | CGG | TCC | GTT | AAC | TCT | G |  
| Rsr II |

30

35

| Hpa I |

| i | y | n | r | v | m | e | s | f | k |  
| 137 | 138 | 139 | 140 | 141 | 142 | 143 | 144 | 145 | 146 |  
AG | ATC | TAT | AAT | CGC | GTT | ATG | GAG | TCG | TTC | AAG |  
| Bgl II |

40

| k | e | g | r | i | g | a | . | . | . |  
| 147 | 148 | 149 | 150 | 151 | 152 | 153 | | | |  
| AAA | GAG | GGT | CGT | ATC | GGC | GCA | TAA | TAG | TGA |

45

50

| GGT | ACC |  
| Kpn I |

55

Table 506: Second Target Downstream  
of Promoters of Selectable Genes

Second Target downstream  
of Pamp that promotes galT.K

5' | CCT | GCG | AAC | CGG | AAT | TGC | CAG | -  
Olig#541 = 3' gga cgc ttg gcc tta acc gtc-  
| StuI | | -35 |

| CTG | GGG | CGC | CCT | CTG | GTA | AGG | TTG | GGA | -  
gac ccc gcg gga gac cat tcc aac cct -  
| -10 |

1024

↓

| TAC | ATA | ATA | CAG | GCA | CTA | TCC | T | A 3' = Olig#542  
atg tat tat gtc cgt gat agg a t tcga 5'  
| Second Target | | Hind3 |

Second Target downstream  
of P<sub>neo</sub> that promotes tet

5' | CTT | CTA | AAT | ACA | TTC | AAA | -  
Olig#543 3' c cgg gaa gat tta tgt aag ttt-  
| ApaI | | -35 |

| TAT | GTA | TCC | GCT | CAT | GAG | ACA | ATA | ACC | CT | -  
ata cat agg cga gta ctg tgt tat tgg ga-  
| -10 |

1024

↓

| TAC | ATA | ATA | CAG | GCA | CTA | TCC | T | CGT 3' = Olig#544  
atg tat tat gtc cgt gat agg a gca gat c 5'  
| Second Target | | XbaI |

Table 507 : Variegation f r selection with  
Second Target

5

10

R|k

| m | k | d | i | w e | g

| 1 | 96 | 97 | 98 | 99 | 100 |

5'-cga|ctg|cgg|TAA|CCT|ATG|AAA|GAT|ATC|TGG|rrA|-

| spacer | BstE II |

15

20

\* \* \* \* \*

M|r K|q G|d M|i S|r K|q

v|g t|p h|r f|l n|k t|p

|101|102|103|104|105|106|

|rkG|mmG|srT|wTk|Ark|mmA|-

25

30

\* \*

M|r P|q H|y F|y

v|g r|l s|p v|d n | l | r | Center of symmetry

|107|108|109|110|111|112|113| ↓ for priming

|rkG|CnG|ymT|kWT|AAC|CTC|AGG|cgt|att|aat|acg|cct|g-3'

| Bsu36I | | primer |

35

k = equimolar T and G                      r = equimolar A and G

w = equimolar T and A                      s = equimolar C and G

m = equimolar A and C                      y = equimolar T and C

40

Approximately  $4 \times 10^6$  DNA and protein sequences.

45

\* indicates sites of one alternative variegation.

50

55

Tabl 508: Protein Ped-6-2  
Selected for Binding to S cond Target

5

| m | k | d | i | w | E |  
| 1 | 96 | 97 | 98 | 99 | 100 |  
| ATG | AAG | GAT | ATT | TGG | GAG |

10

| R | Q | G | M | R | T |  
| 101 | 102 | 103 | 104 | 105 | 106 |  
| AGG | CAG | GGT | ATG | AGG | ACA |

15

| M | P | Y | F | n | l | r | w | p | r |  
| 107 | 108 | 109 | 110 | 111 | 112 | 113 | 114 | 115 | 116 |  
| ATG | CCT | TAC | TTT | AAC | CTC | AGG | TGG | CCC | CGG | G |  
| Bsu36I | | Xma I |

20

| e | v | l | d | l | v | r | k | v | a |  
| 117 | 118 | 119 | 120 | 121 | 122 | 123 | 124 | 125 | 126 |  
| AG | GTC | CTT | GAT | CTT | GTT | CGC | AAG | GTT | GCT |  
| PpuM I |

25

| e | e | n | g | r | s | v | n | s | e |  
| 127 | 128 | 129 | 130 | 131 | 132 | 133 | 134 | 135 | 136 |  
| GAG | GAA | AAC | GGT | CGG | TCC | GTT | AAC | TCT | G |  
| Rsr II |

30

35

| Hpa I |

| i | y | n | r | v | m | e | s | f | k |  
| 137 | 138 | 139 | 140 | 141 | 142 | 143 | 144 | 145 | 146 |  
AG | ATC | TAT | AAT | CGC | GTT | ATG | GAG | TCG | TTC | AAG |  
| Bgl II |

40

| k | e | g | r | i | g | a | . | . | . |  
| 147 | 148 | 149 | 150 | 151 | 152 | 153 | | | |  
| AAA | GAG | GGT | CGT | ATC | GGC | GCA | TAA | TAG | TGA |

45

50

| GGT | ACC |  
| Kpn I |

55

Table 509: Protein Ped-6-2-5  
Selected for Binding to Third Target

5

| m | R | D | V | W | H |  
| 1 | 96 | 97 | 98 | 99 | 100 |  
| ATG | AGG | GAT | GTT | TGG | CAT |

10

| V | R | N | I | T | R |  
| 101 | 102 | 103 | 104 | 105 | 106 |  
| GTG | CGG | AAT | ATT | ACG | AGA |

15

| V | R | H | L | n | l | r | w | p | r |  
| 107 | 108 | 109 | 110 | 111 | 112 | 113 | 114 | 115 | 116 |  
| GTG | CGT | CAC | TTG | AAC | CTC | AGG | TGG | CCC | CGG | G |  
| Bsu36I | | Xma I |

20

| e | v | l | d | l | v | r | k | v | a |  
| 117 | 118 | 119 | 120 | 121 | 122 | 123 | 124 | 125 | 126 |  
| AG | GTC | CTT | GAT | CTT | GTT | CGC | AAG | GTT | GCT |  
| PvuM I |

25

| e | e | n | g | r | s | v | n | s | e |  
| 127 | 128 | 129 | 130 | 131 | 132 | 133 | 134 | 135 | 136 |  
| GAG | GAA | AAC | GGT | CGG | TCC | GTT | AAC | TCT | G |  
| Rsr II |

30

35

| Hpa I |

| i | y | n | r | v | n | e | s | f | k |  
| 137 | 138 | 139 | 140 | 141 | 142 | 143 | 144 | 145 | 146 |  
AG | ATC | TAT | AAT | CGC | GTT | ATG | GAG | TCG | TTC | AAG |  
| Bgl II |

40

| k | e | g | r | i | g | a | . | . | . |  
| 147 | 148 | 149 | 150 | 151 | 152 | 153 | | | |  
| AAA | GAG | GGT | CGT | ATC | GGC | GCA | TAA | TAG | TGA |

45

| GGT | ACC |  
| Kpn I |

50

55

Table 510: Variegation f Ped-6-2  
for binding to Third Target

5  
10  
K|r R|h  
m t|i D|n I|v W|l p|l  
| 1 | 96| 97| 98| 99|100|  
|cct|cag|cGG|TAA|CCT|ATG|AnA|rAT|rTT|TmG|CnT|  
| spacer | | BstE II |

15  
20  
G|s  
V|d Q|r n|d M|i R|t T|r  
|101|102|103|104|105|106|  
|GwT|CrG|rrT|ATr|AsG|AsA|

25  
30  
M|k F|c  
e|v P|r Y|h l|w n | l | r | Center of symmetry  
|107|108|109|110|111|112|113| ↓ for priming  
|rwG|CsT|yAC|Tkk|AAC|CTC|AGG|cgt|att|aat|acg|cct|g -3'  
| Bsu36I | olig#506 ↑

35  
Self priming for extension with Klenow

40  
5'-..CTC|AGG|cgt|att|\  
3'- g|tcc|gca|taa|/



Table 511: Variegati n for  
Selection with F urth Target

5

10

|        |     |     |         |     |     |     |     |
|--------|-----|-----|---------|-----|-----|-----|-----|
| m      | X   | X   | X       | X   | X   | X   | X   |
| 1      | 90  | 91  | 92      | 93  | 94  | 95  |     |
| cct    | cag | cGG | TAA     | CCT | ATG | fzk | fzk |
|        |     |     |         |     |     | fzk | fzk |
| spacer |     |     | BstE II |     |     |     |     |

15

|         |     |     |     |     |
|---------|-----|-----|-----|-----|
| R       | D   | V   | W   | H   |
| 96      | 97  | 98  | 99  | 100 |
| CGG     | GAC | GTG | TGG | CAC |
| overlap |     |     |     |     |

20

25

|     |     |     |     |     |     |
|-----|-----|-----|-----|-----|-----|
| V   | R   | N   | I   | T   | R   |
| 101 | 102 | 103 | 104 | 105 | 106 |
| GTG | CGG | AAT | ATT | ACG | CGA |

30

|               |     |     |     |     |     |     |
|---------------|-----|-----|-----|-----|-----|-----|
| V             | R   | H   | L   | n   | l   | r   |
| 107           | 108 | 109 | 110 | 111 | 112 | 113 |
| GTG           | CGT | CAC | CTT | AAC | CTC | AGG |
| cgt           | cac | ggc |     |     |     |     |
| Bsu36I spacer |     |     |     |     |     |     |

35

f = (0.26 T, 0.18 C, 0.26 A, 0.30 G)

z = (0.22 T, 0.16 C, 0.40 A, 0.22 G)

40

k = equimolar T and G

There are  $(2^5)^6 = 2^{30} = 10^9$  DNA sequences.

45

There are  $20^6 = 6.4 \times 10^7$  protein sequences.

50

55

Table 512: Prot in Pad-6-2-5-2  
Selected for Binding to Fourth Target

5

| m | R | T | G | F | C | Q |  
| 1 | 90 | 91 | 92 | 93 | 94 | 95 |  
| ATG | CGT | ACG | GGG | TTT | TGT | CAG |

10

| R | D | V | W | H |  
| 96 | 97 | 98 | 99 | 100 |  
| CGG | GAT | GTT | TGG | CAC |

15

| V | R | N | I | T | R |  
| 101 | 102 | 103 | 104 | 105 | 106 |  
| GTG | CGG | AAT | ATT | ACG | CGA |

20

| V | R | H | L | n | l | r | w | p | r |  
| 107 | 108 | 109 | 110 | 111 | 112 | 113 | 114 | 115 | 116 |  
| GTG | CGT | CAC | CTT | AAC | CTC | AGG | TGG | CCC | CGG | G |  
| Bsu36I | | Xma I |

25

30

| e | v | l | d | l | v | r | k | v | a |  
| 117 | 118 | 119 | 120 | 121 | 122 | 123 | 124 | 125 | 126 |  
| AG | GTC | CTT | GAT | CTT | GTT | CGC | AAG | GTT | GCT |  
| PvuM I |

35

| e | e | n | g | r | s | v | n | s | e |  
| 127 | 128 | 129 | 130 | 131 | 132 | 133 | 134 | 135 | 136 |  
| GAG | GAA | AAC | GGT | CGG | TCC | GTT | AAC | TCT | G |  
| Rsr II |

40

45

| Hpa I |

| i | y | n | r | v | m | e | s | f | k |  
| 137 | 138 | 139 | 140 | 141 | 142 | 143 | 144 | 145 | 146 |  
AG | ATC | TAT | AAT | CGC | GTT | ATG | GAG | TCG | TTC | AAG |  
| Bgl II |

50

55

Table 512, continued

5           | k | e | g | r | i | g | a | . | . | . |  
           |147|148|149|150|151|152|153|   |   |   |  
           |AAA|GAG|GGT|CGT|ATC|GGC|GCA|TAA|TAG|TGA|

10

          |GGT|ACC|  
           | Kpn I |

15

20

25

30

35

40

45

50

55

Table 513: Variegati n  
of Length of Ped-6-2-5-2

5  
10  
N|k V|g M|k  
i y|. y|. r|. l|. l|. |  
|137| 138| 139 | 140| 141 | 142 |  
5' = cgacctagcAG|ATC|TAW<sub>1</sub>|w<sub>3</sub>Aw<sub>4</sub>|y<sub>1</sub>GA|k<sub>1</sub>k<sub>2</sub>A|w<sub>3</sub>w<sub>4</sub>G|  
| spacer | Bgl II |

15  
F|c  
e|. s|. l|. k|. |  
| 143| 144| 145 | 146|  
20 |k<sub>3</sub>AG|Tm<sub>1</sub>G|Tk<sub>2</sub>m<sub>2</sub>|w<sub>2</sub>AG|-

25  
I|k  
k|. e|. g|. r|. l|. g|. | . | . |  
| 147| 148| 149| 150| 151 | 152|153| |  
30 |w<sub>2</sub>AA|k<sub>3</sub>AG|k<sub>3</sub>GA|y<sub>1</sub>GA|w<sub>3</sub>w<sub>4</sub>A|k<sub>3</sub>GA|TAG|TGA|-

|GGT|ACC|t- 3'  
| Kpn I |

35  
w<sub>1</sub> = 0.65 T and 0.35 A y<sub>1</sub> = 0.65 C and 0.35 T  
k<sub>1</sub> = 0.42 G and 0.58 T k<sub>2</sub> = 0.42 T and 0.58 G  
40 k<sub>3</sub> = 0.65 G and 0.35 T m<sub>1</sub> = 0.65 C and 0.35 A  
m<sub>2</sub> = 0.42 C and 0.58 A w<sub>2</sub> = 0.65 A and 0.35 T  
w<sub>3</sub> = 0.42 A and 0.58 T w<sub>4</sub> = 0.42 T and 0.58 A

45 Each variegated residue produces about 35% stop codons.

50 Because  $(0.65)^{15} = 0.003$ , only 0.3 % of variegated genes  
encode a protein shortened by one residue.

Table for Example 6

Table 600: Third finger domain of kr -tgs- P22 arc

5

10

15

20

25

30

35

40

45

50

55

```

      *
      | m | e | k | p | y | h |
      | 1 | 2 | 3 | 4 | 5 | 6 |
|AGG|AGG|TAA|CCT|ATG|GAG|AAA|CCG|TAT|CAC|-
  |BstE II|

      *           *   *
      | c | s | h | c | d | r | q | E | v | q |
      | 7 | 8 | 9 | 10| 11| 12| 13| 14| 15| 16|
|TGC|TCA|CAC|TGT|GAT|CGT|CAG|TTT|GTC|CAA|-
  |Dra III|

      *   *           *   *   *
      | v | a | n | l | r | r | H | l | r | v | H |
      | 17| 18| 19| 20| 21| 22| 23| 24| 25| 26| 27|
|GTG|GCC|AAC|TTA|AGA|CGT|CAT|CTA|CGC|GTG|CAC|-
  |Bal I| |Afl II|Aat II| |Mlu I|
                                |ApaL I|

      |<- linker >|<---- P22 arc
      *   *   *   *   *   *   *   *
      | t | g | t | g | s | m | k | g | m |
      | 28| 29| 30| 31| 32|101|102|103|104|
|ACT|GGT|ACC|GGG|TCT|ATG|AAA|GGC|ATG|
  |Kpn I|

      *   *   *   *   *   *
      | s | k | m | p | q | f | n | l | r | w |
      |105|106|107|108|109|110|111|112|113|114|
|TCT|AAG|ATG|CCG|CAA|TTC|AAC|CTT|AGG|TGG|
  |Bsu36I|

```

Tabl 600, continued

5 | p | r | e | v | l | d | l | v | r | k |  
 | 115 | 116 | 117 | 118 | 119 | 120 | 121 | 122 | 123 | 124 |  
 | CCC | CGG | GAG | GTC | CTT | GAT | TTG | GTT | CGC | AAA |  
 10 | Ava I | PvuII |  
 | Xma I |  
 15 | v | a | e | e | n | g | r | s | v | n | s |  
 | 125 | 126 | 127 | 128 | 129 | 130 | 131 | 132 | 133 | 134 | 135 |  
 | GTC | GCT | GAA | GAG | AAT | GGC | CGG | TCC | GTG | AAT | TCT |  
 20 | Ksp 632 | Rsr II | EcoR I |  
 | e | i | y | n | r | v | m | e | s |  
 25 | 136 | 137 | 138 | 139 | 140 | 141 | 142 | 143 | 144 |  
 | GAG | ATC | TAT | AAT | CGT | GTT | ATG | GAA | AGC |  
 | Bgl II |  
 30 | f | k | k | e | g | r | i | g | a | . |  
 | 145 | 146 | 147 | 148 | 149 | 150 | 151 | 152 | 153 | 154 |  
 | TTC | AAG | AAG | GAA | GGT | CGC | ATT | GGT | GCA | TAA |  
 35 | . | . |  
 | 155 | 156 |  
 40 | TAG | TGA | GGA | TTC |  
 | HindIII |

45 \* indicates residues of zinc finger domain thought to  
 contact DNA in model of Gibson et al.

50 \* indicates residues of zinc finger domain, linker, and  
 Arc that may influence DNA binding.

55

## CITATIONS:

[0407]

- ADHY82 Adhya, S, and M Gottesman.  
Cell (1982) 29:939-944.
- AGGA88 Aggarwal, AK, DW Rodgers, M Drott, M Ptashne, SC Harrison. Science (1988), 242:899-907.
- AHME84 Ahmed, A. Gene (1984) 28:37-43.
- ANDE81 Anderson, WF, DH Ohlendorf, Y Takeda, BW Matthews. Nature (1981), 290:754-758.
- 10 ANDE87 Anderson, JE, M Ptashne, and SC Harrison.  
Nature (1987), 326:846-852.
- ANDE88 Anderson, WF, M Cygler, RP Braun, and JS Lee.  
Bioessays (1988), 8:69-74
- AUSU87 Ausubel, FM, R Brent, RE Kingston, DD Moore, JG Seidman, JA Smith, K Struhl.  
15 Current Protocols in Molecular Biology. 1987, John Wiley and Sons.
- BASH87 Bash, PA, UC Singh, R Langridge, and PA Kollman. Science (1987), 236(4801):564-8.
- BASS88 Bass, S, V Sorrells, and P Youderian.  
Science (1988), 242:240-245.
- BECK88 Becker, MM, D Lesser, M Kurpiewski, A Baranger, and L Jen-Jacobson. Proc Natl Acad Sci USA (Sept.  
20 1988), 85:6247-6251.
- BENS86 Benson, N, P Sugiono, S Bass, LV Mendelman, and P Youderian. Genetics (1986), 114:1-14.
- BENS88 Benson, N, P Sugiono, and P Youderian.  
Genetics (1988), 188:21-29.
- BERG88a Berg, JM.  
25 Proc Natl Acad Sci USA (1988), 85:99-102.
- BOCH80 Bochner, BR, H Huang, GL Schieven, and BN Ames.  
J Bacteriol (1980), 143:926-933.
- BOHN88 Bohnlein, E, JW Lowenthal, M Siekevitz, DW Ballard, BR Franza, and WC Greene.  
Cell (1988), 53:827-836.
- 30 BOTS85 Botstein, D, and D Shortle.  
Science (1985), 229:1193-1201.
- BREN84 Brent, R, and M Ptashne.  
Nature (1984), 312:612-615.
- BROS82 Brosius, J, RL Cate, and AP Perlmutter.  
35 J Biol Chem (1982), 257:9205-9210.
- BROS84 Brosius, J. Gene (1984), 27:151-160.
- BROW87 Brown, M, J Figge, U Hansen, C Wright, KT Jeang, G Khoury, DM Livingston, and TM Roberts.  
Cell (1987), 49:603-612.
- BRUN87 Brunelle, A, and RF Schleif.  
40 Proc Natl Acad Sci USA (1987), 84:6673-6676.
- BUSH85 Bushman, FD, JE Anderson, SC Harrison, and M Ptashne. Nature (1985), 316:651-653.
- BUSH88 Bushman, FD, and M Ptashne.  
Cell (1988), 54:191-197.
- BUTT63 Buttin, G. J Mol Biol (1963), 7:164-182.
- 45 CARU85 Caruthers, MH. Science (1985), 230:281-285.
- CARU87 Caruthers, MH, P Gottlieb, L Bracco, and L Cummins. Protein Structure, Folding, and Design (1987), 2:9-24.
- CHAD71 Chadwick, P, V Pirrotta, R Steinberg, N Hopkins, and M Ptashne. Cold Spring Harb Symp Quant Biol (1971), 35:283-294.
- 50 CHEN88 Chen, W, and K Struhl.  
Proc Natl Acad Sci USA (1988), 85:2691-2695.
- CHOU78a Chou, PY, and GD Fasman.  
Adv Enzymol (1978), 47:45-148.
- CHOU78b Chou, PY, and GD Fasman.  
55 Annu Rev Biochem (1978), 47:251-76.
- CRAI85 Craik, CS, C Largman, T Flecher, S Rocznia, PJ Barr, R Fletterick, and WJ Rutter.  
Science (1985) 228:291-297.
- DAVI80 Davis, RW, D Botstein, and JR Roth. Advanced Bacterial Genetics. Cold Spring Harbor Laboratory Press.

1980.

- DAYR86 Dayringer, H, A Tramantano, and R Fletterick. p.5-8 in Computer Graphics and Molecular Modeling, Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, 1986
- DAYT86 Dayton, AI, JG Sodroski, CA Rosen, WC Goh, and WA Haseltine. *Cell* (1986), 44:941-947.
- 5 DEFE86 DeFeyter, RC, BE Davidson, and J Pittard.  
*J Bacteriol* (1986), 165:233-239.
- DELO87 De Lorenzo, V, S Wee, M Herrero, and JB Neilands. *J Bacteriol* (1987) 165:2624-2630.
- DILL87 Dill, KA.  
*Protein Engineering* (1987), 1:369-371.
- 10 EBRI84 Ebright, RH, P Cossart, B Gicquel-Sanzey, and J Beckwith.  
*Proc Natl Acad Sci USA* (1984), 81:7274-7278.
- EISE85 Eisenbeis, SJ, MS Nasoff, SA Noble, LP Bracco, DR Dodds, and MH Caruthers.  
*Proc Natl Acad Sci USA* (1985), 82:1084-1088.
- EISE86a Eisenberg, D, W Wilcox, and AD McLachlan.  
15 *J cell Biochem* (1986), 31:11-17.
- EISE86b Eisenberg, D, and AD McLachlan.  
*Nature* (1986), 319:199-203.
- ELIA85 Eliason, JL, MA Weiss, and M Ptashne.  
*Proc Natl Acad Sci USA* (1985), 82:2339-2343.
- 20 ELLE89a Elledge, SJ, and RW Davis.  
*Genes & Development* (1989) 3:185-197.
- ELLE89b Elledge, SJ, P Sugiono, L Guarente, and RW Davis.  
*Proc Natl Acad Sci USA* (1989) 86:3689-93.
- EVAN88 Evans, RM, and SM Hollenberg.  
25 *Cell* (1988), 52:1-3.
- FAIR86 Fairall, L, D Rhodes, and A Klug.  
*J Mol Biol* (1986), 192:577-591.
- FEIN86 Feinberg, MB, RF Jarrett, A Aldovini, RC Gallo, and F Wong-Staal. *Cell* (1986), 46:807-817.
- FIGG88 Figge, J, C Wright, CJ Collins, TM Roberts, and DM Livingston. *Cell* (1988), 52:713-722.
- 30 FILU85 Filutowicz, M, G Davis, A Greener, and DR Helinski. *Nucl Acids Res* (1985), 13:103-114.
- FOXK88 Fox, KR.  
*Biochem & Biophys Res Comm* (1988), 155(2):779-85.
- FRAN88 Frankel, AD, and CO Pabo. *Cell* (1988), 53:675.
- FRIE81 Fried, MG, and DM Crothers.  
35 *Nucl Acids Res* (1981), 9:6505-6525.
- GART88 Gartenberg, MR, and DM Crothers.  
*Nature* (1988), 333:824-829.
- GIBS88 Gibson, TJ, JPM Postma, RS Brown, and P Argos. *Protein Engineering* (1988), 2:209-218.
- GUAR82 Guarente, L, JS Nye, A Hochschild, and M Ptashne.  
40 *Proc Natl Acad Sci USA* (1982), 79:2236-2239.
- HANA85 Hanahan, D. (In:DNA Cloning Volume I. 1985 IRL Press) pp109-135.
- HARR88 Harrison, SC, JE, anderson, GB Koudelka, A Mondragon, S Subbiah, RP Wharton, C Wolberger, and M Ptashne. *Biophys Chem* (1988), 29:31-37.
- HAWL83 Hawley, DK, and WR McClure.  
45 *Nucl Acids Res* (1983), 11(8):2237-2255.
- HECH83 Hecht, MH, HCM Nelson, and RT Sauer.  
*Proc Natl Acad Sci USA* (1983), 80:2676-2680.
- HECH84 Hecht, MH, JM Sturtevant, and RT Sauer.  
*Proc Natl Acad Sci USA* (1984), 81:5685-5689.
- 50 HECH85a Hecht, MH, and RT Sauer.  
*J Mol Biol* (1985), 186:53-63.
- HECH85b Hecht, MH, KM Hehir, HCM Nelson, JM Sturtevant, and RT Sauer.  
*J Cellular Biochem* (1985), 29:217-224.
- HOCH83 Hochschild, A, N Irwin, and M Ptashne.  
55 *Cell* (1983), 32:319-325.
- HOCH86a Hochschild, A, and M Ptashne.  
*Cell* (1986), 44:925-933.
- HOCH86b Hochschild, A, J Douhan III, and M Ptashne.



- Cell (1986), 47:807-816.
- HOGA87 Hogan, ME, and RH Austin.  
Nature (1987), 329:263-266.
- HOLL88 Hollis, M, D Valenzuela, D Pioli, R Wharton, and M Ptashne.  
Proc Natl Acad Sci USA (1988), 85:5834-5838.
- 5 HOOP87 Hoopes, BC, and WR McClure.  
Volume 2, Chapter 75, pp1231-1240, Escherichia coli, and Salmonella typhimurium: Cellular, and Molecular Biology, Neidhardt, FC, Editor-in-Chief. Amer. Soc. for Microbiology, Washington, DC, 1987.
- HUMC87 Hu, MCT, and N Davidson.  
10 Cell (1987), 48:555-566.
- HUMC88 Hu, MCT, and N Davidson.  
Gene (1988), 62:301-313.
- INOUE86 Inouye, M, and R Sarma, Editors. Protein Engineering: Applications in Science, Medicine, and Industry, Academic Press, New York, 1986.
- 15 JENJ86 Jen-Jacobson, L, D Lesser, and M Kurpiewski. Cell (1986), 45:619-629.
- JOHN79 Johnson, AD, BJ Meyer, and M Ptashne.  
Proc Nat Acad Sci USA (1979), 10:5061-5065.
- JOHN80 Johnson, AD, CO Pabo, and RT Sauer.  
Meth Enzymol (1980), 65: 839-856.
- 20 JOHN86 Johnson, HL, MR Gartenberg, and DM Crothers.  
Cell (1986), 47:995-1005.
- JONE85 Jones, TA. Methods Enzymol (1985), 115:157-71.
- JONE87 Jones, KA, JT Kadonaga, PJ Rosenfeld, TJ Kelly, and R Tjian. Cell (1987), 48:79-89.
- JORD85 Jordan, SR, CO Pabo, AK Vershon, and RT Sauer.  
25 J Mol Biol (1985), 185:445-446.
- KADO86 Kadonaga, JT, and R Tjian.  
Proc Natl Acad Sci USA (1986), 83:5889-5893.
- KARN84 Karn, J, HWD Matthes, MJ Gait, and S Brenner.  
Gene (1984), 32:217-224.
- 30 KELL85 Kelley, RL, and C Yanofsky.  
Proc Natl Acad Sci USA 82:483-487.
- KIMJ87 Kim, JG, Y Takeda, BW Matthews, and WF Anderson.  
J Mol Biol (1987), 196:149-158.
- KLEN70 Klenow, H, and I Henningsen.  
35 Proc Natl Acad Sci USA (1970), 65:168-175.
- KNIG88 Knight, KL, and RT Sauer.  
Biochem (1988), 27:2088-2094.
- KOUD87 Koukelka, GB, SC Harrison, and M Ptashne.  
Nature (1987), 326:886-888.
- 40 KOUD88 Koudelka, GB, P Harbury, SC Harrison, and M Ptashne.  
Proc Natl Acad Sci USA (1988), 85:4633-4637.
- KRAU86 Krause, HM, and NP Higgins.  
J Biol Chem (1986), 261:3744-3752.
- LATH85 Lathe, R. J Mol Biol (1985), 183:1-12.
- 45 LEGE85 Legerski, RJ, and DL Robberson.  
J Mol Biol (1985), 181:297-312.
- LEIG87 Leighton, P, and P Lu.  
Biochem (1987), 26:7262-7271.
- LEWI83 Lewis, M, A Jeffrey, J Wang, R Ladner, M Ptashne, and CO Pabo. Cold Spring Harbor Symp Quant Biol  
50 (1983), 47:435-440.
- LIPM85 Lipman, DJ, and WR Pearson.  
Science (1985), 227:1435-1441.
- MALO81 Maloy, SR, and WD Nunn.  
J Bacteriol (1981), 145:1110-1111.
- 55 MANI82 Maniatis, T, EF Fritsch, and J Sambrook.  
Molecular Cloning, Cold Spring Harbor Laboratory, 1982.
- MANI87 Maniatis, T, S Goodbourn, and JA Fischer. Science (1987), 236:1237-1245.
- MATT88 Matthews, BW. Nature (1988), 335:294-295.

- MAUR80 Maurer, R, BJ Meyer, and M Ptashne.  
J Mol Biol (1980), 139:147-161.
- MAXA77 Maxam, A, and W Gilbert.  
Proc Natl Acad Sci USA (1977), 74:560-564.
- 5 MAXA80 Maxam, A, and W Gilbert.  
Meth Enzymol (1980), 65:499-599.
- MCCL86 McClarin, CA Frederick, B-C Wang, P Greene, HW Boyer, J Grable, and JM Rosenberg.  
Science (1986), 234:1526-41.
- MCKA81 McKay, DB, and TA Steitz.  
10 Nature (1981), 290:744-749.
- MCKA82 McKay, DB, IT Weber, and TA Steitz.  
J Biol Chem (1982), 257:9518-9524.
- MILL72 Miller, JH. Experiments in Molecular Genetics. Cold Spring Harbor Laboratory, Cold Spring Harbor, NY.  
1972.
- 15 MILL85 Miller, AM, VL Mackay, and KA Nasmyth.  
Nature (1985), 314:598-603.
- MULL84 Mulligan, ME, DK Hawley, R Entriken, and WR McClure.  
Nucleic Acids Research (1984), 12:789-800.
- NEID87b Neidle S, LH Pearl, and JV Skelly.  
20 Biochem J (1987), 243:1-13.
- NELS83 Nelson, HCM, MH Hecht, and RT Sauer.  
Cold Spring Harbor Symp Quant Biol (1983), 47:441-449.
- NELS85 Nelson, HCM, and RT Sauer.  
Cell (1985), 42:549-558.
- 25 NELS86 Nelson, HCM, and RT Sauer.  
J Mol Biol (1986), 192:27-38.
- NIKA61 Nikaido, H.  
Biochem Biophys Acta (1961), 48:460-469.
- OHLE83 Ohlendorf, DH, WF Anderson, M Lewis, CO Pabo, and W Matthews. J Mol Biol (1983), 169:757-769.
- 30 OHLE85 Ohlendorf, DH, and JB Matthew.  
Adv Biophys (1985), 20:137-51.
- OLIP86 Oliphant, AR, AL Nussbaum, and K Struhl.  
Gene (1986), 44:177-183.
- OLIP87 Oliphant, AR, and K Struhl.  
35 Methods in Enzymology 155 (1987):568-582. Editor Wu, R; Academic Press, New York, 1987.
- OTWI88 Otwinoski, Z, RW Schevitz, R-G Zhang, CL Lawson, A Joachimiak, RQ Marmorstein, BF Luisi, and PB  
Sigler. Nature (1988), 335:321-329.
- PABO79 Pabo, CO, RT Sauer, JM Sturtevant, and  
M Ptashne. Proc Natl Acad Sci USA (1979), 76:1608-1612.
- 40 PABO82a Pabo, CO, W Krovatin, A Jeffrey, and RT Sauer.  
Nature (1982), 298:441-443.
- PABO82b Pabo, CO, and M Lewis.  
Nature (1982), 298:443-447.
- PABO84 Pabo, CO, and RT Sauer.  
45 Ann Rev Biochem (1984), 53:293-321.
- PAKU86 Pakula, AA, VB Young, and RT Sauer.  
Proc Natl Acad Sci USA (1986), 83:8829-8833.
- PARR88 Parraga, G, SJ Horvath, A Eisen, WE Taylor, L Hood, T Young, RE Klevit.  
Science (1988), 241:1489-1492.
- 50 POLA88 Polayes, DA, PW Rice, MM Garner, and JE Dahlberg. J Bacteriol (1988), 170:3110-3114.
- POTE80 Poteete, AR, M Ptashne, M Ballivet, and H Eisen. J Mol Biol (1980), 137:81-91.
- POTE82 Poteete, AR, and M Ptashne.  
J Mol Biol (1982), 157:21-48.
- PTAS80 Ptashne, M, A Jeffrey, AD Johnson, R Maurer, BJ Meyer, CO Pabo, TM Roberts, and RT Sauer.  
55 Cell (1980), 19:1-11.
- PTAS86 Ptashne, M. A Genetic Switch: Gene Control, and Phage λ. 1986, Cell Press, and Blackwell Scientific Pub-  
lications.
- RAOS87 Rao, SN, UC Singh, PA Bash, and PA Kollman.

- Nature (1987), 328(6130) 551-4.
- RATN85 Ratner, L, W Haseltine, KR Patarca, KJ Livak, B Starcich, SF Josephs, ER Doran, JA Rafalski, EA Whitehorn, K Baumeister, L Ivanoff, SLR Petteway Jr, ML Pearson, JA Lautenberger, TS Papas, J Ghayeb, NT Chang, RC Gallo, and F Wong-Staal. Nature (1985), 313:277-284.
- 5 REID88 Reidhaar-Olson, JF, and RT Sauer. Science (1988), 241:53-57.
- RENY88 Ren, YL, S Garges, S Adhya, and JS Krakow. Mol Microbiol (1987), 1:53-58.
- RICH86 Richards, JH. Nature (1986), 323:187.
- 10 RICH88 Richet, E, P Abcarian, and HA Nash. Cell (1988), 52:9-17.
- RIGB77 Rigby, PW, M Dieckmann, C Rhodes, and P Berg. J Mol Biol (1977), 113:237-251.
- RIGG70 Riggs, AD, H Suzuki, and S Bourgeois. J Mol Biol (1970), 48:67-83.
- 15 ROBE86 Roberts, S, and AR Rees. Protein Engineering (1986), 1:59-65.
- ROSE79 Rosenberg, M, and D Court. Ann Rev Genet (1979), 13:319-353.
- 20 ROSE85 Rose, GD, AR Geselowitz, GJ Lesser, RH Lee, and MH Zehfus. Science (1985), 229:834-838.
- ROSE86 Rosenberg, UB, C Schroeder, A Preiss, A Kienlin, S Cote, I Riede, and H Jaeckle. Nature (1986), 319:336-339.
- SAAG88 Saag, MS, BH Hahn, J Gibbons, Y Li, ES Parks, WP Parks, and GM Shaw. Nature (1988), 334:440-444.
- SADL83 Sadler, JR, H Sasmor, and JL Betz. Proc Natl Acad Sci USA (1983), 80:6785-6789.
- 25 SAEN83 Saenger, W. Principles of Nucleic Acid Structure. Springer Verlag, New York, 1983.
- SAUE79 Sauer, RT, CO Pabo, BJ Meyer, M Ptashne, and KC Backman. Nature (1979), 279:396-400.
- SAUE82 Sauer, RT, RR Yocum, RF Doolittle, M Lewis, and CO Pabo. Nature (1982), 298:447-451.
- SAUE86 Sauer, RT, K Hehir, RS Stearman, MA Weiss, A Jeitler-Neilsson, EG Suchanek, and CO Pabo. Biochem (1986), 25:5992-5998.
- 30 SCHL88 Schleif, R. Science (1988), 241:1182-87.
- SIMO84 Simons, A, D Tils, B von Wilcken-Bergmann, and B Muller-Hill. Biochem. (1984) 81:1624-1628.
- SIMO88 Simon, MC, TM Fisch, BJ Benecke, JR Nevins, and N Heintz. Cell (1988), 52:723-729.
- SMIT87 Smith, DI, W Golembieski, JD Gilbert, L Kizyma, and OJ Miller. Nucl Acid Res (1987), 15(3):1173-84.
- 35 SOUT75 Southern, E. J. Mol. Biol. (1975) 98:503.
- SPIR88 Spiro, S, and JR Guest. Mol Microbiol (1987), 1:53-58.
- STRU87 Struhl, K. Cell (1987), 49:295-297.
- 40 TAKE77 Takeda, Y, A Folkmanis, and H Echols. J Biol Chem (1977), 252:6177-6183.
- TAKE83 Takeda, Y, H Ohlendorf, WF, anderson, and BW Matthews. Science (1983), 221:1020-1026.
- TAKE85 Takeda, Y, DH Ohlendorf, WF Anderson, and BW Matthews. In: Biological Macromolecules, and Assemblies: Volume 2-Nucleic Acids, and Interactive Proteins. 1985, John Wiley, and Sons, Inc. pp234-263.
- 45 TAKE86 Takeda, Y, JG Kim, CG Caday, E Steers Jr., DH Ohlendorf, WF, anderson, and BW Matthews. J Biol Chem (1986), 261:8608-8616.
- THER88 Theriault, NY, JB Carter, and SP Pulaski. Biotechniques (1988), 6:470-474.
- ULAN87 Ulanovsky, LE, and EN Trifonov. Nature (1987), 326:720-722.
- 50 ULME83 Ulmer, KM. Science (1983), 219(4585):666-71.
- VERS85a Vershon, AK, P Youderian, MA Weiss, MM Susskind, and RT Sauer. In: Sequence Specificity in Transcription, and Translation. 1985, Alan R. Liss, Inc. pp209-218.
- VERS85b Vershon, AK, P Youderian, MM Susskind, and RT Sauer. J Biol Chem (1985), 260:12124-12129.
- 55 VERS86a Vershon, AK, K Blackmer, and RT Sauer. In: Protein Engineering. Applications, in Science, Medicine, and Industry. 1986, Academic Press, Inc. p243-256.
- VERS86b Vershon, AK, JU Bowie, TM Karplus, and RT Sauer. In: Proteins: Structure, Function, and Genetics. 1986,

Alan R. Liss, Inc. pp302-311.

VERS87a Vershon, AK, SM Liao, WR McClure, and RT Sauer.  
J Mol Biol (1987), 195:311-322.

VERS87b Vershon, AK, SM Liao, WR McClure, and RT Sauer.  
J Mol Biol (1987), 195:323-331.

VINO87 Vinopal, RT. In: Escherichia coli, and Salmonella typhimurium: Cellular, and Molecular Biology. 1987, American Society for Microbiology.

WARD86 Ward, WH, DH Jones, and AR Fersht.  
J Biol Chem (1986), 261:9576-8.

10 WEBE87a Weber, IT, GL Giliiland, JG Harman, and A Peterkofsky.  
J Biol Chem (1987), 262:5630-5636.

WEIS87a Weiss, MA, M Karplus, and RT Sauer.  
Biochem (1987), 26:890-897.

WEIS87b Weiss, MA, CO Pabo, M Karplus, and RT Sauer.  
Biochem (1987), 26:897-904.

15 WEIS87c Weiss, MA, M Karplus, and RT Sauer.  
J Biomol Struct Dynam (1987), 5:539-556.

WHAR84 Wharton, RP, EL Brown, and M Ptashne.  
Cell (1984), 38:361-369.

20 WHAR85a Wharton, RP. Ph. D. Thesis, Harvard University.

WHAR85b Wharton, RP, and M Ptashne.  
Nature (1985), 316:601-605.

WHAR87 Wharton, RP, and M Ptashne.  
Nature (1987), 325:888-891.

25 WOLB88 Wolberger, C, Y Dong, M Ptashne, and SC Harrison. Nature (1988), 335:789-795.

YANO87 Yanofsky, SD, R Love, JA McClarin, JA Rosenberg, and HW Boyer. Proteins (1987), 2:273-282.

YOU83 Youderian, P, A Vershon, S Bouvier, RT Sauer, and MM Susskind. Cell (1983), 35:777-783.

ZOLL84 Zoller, MJ, and M Smith. DNA (1984), 3:479-488.

### 30 Claims

1. A selection vector for selecting recipient cells transformed by such vector that express a protein or polypeptide that binds specifically to a predetermined target DNA sequence borne by said vector, which comprises a first and a second operon, each comprising at least one expressible gene, the genes of said first and second operon being different, a copy of the target DNA sequence being included in each operon and positioned therein so that the recipient cells enjoy a selective advantage, other than resistance to lytic growth of phage, if they express a protein or polypeptide which binds to said copies of the target DNA sequence.

2. The vector of claim 1 wherein at least one operon comprises a selectable beneficial gene, an occludible promoter operably linked to said beneficial gene and directing its transcription, an occluding promoter occluding transcription from said occludible promoter of said beneficial gene, and a copy of the target DNA sequence positioned so that the binding of said protein or polypeptide to said copy represses said occluding promoter and thereby facilitates transcription of said beneficial gene.

3. The selection vector of either of claims 1 or 2 which comprises:

a) a first operon, which operon comprises:

i) a first binding marker gene or genes,

ii) a first promoter controlling expression of said binding marker gene or genes, and

iii) a first copy of the target DNA sequence, where said target DNA sequence interferes substantially with expression of the first gene(s) if a protein expressed by the recipient cell binds to the target DNA sequence,

b) a second operon, which operon comprises:

i) a second binding marker gene or genes,

ii) a second promoter controlling expression of said second binding marker gene or genes, and

- iii) a second copy of the target DNA sequence,  
where said target DNA sequence interferes substantially with expression of the second gene(s) if a protein expressed by the recipient cell binds to the target DNA sequence,  
where the binding marker genes of said first and second operons are different, and where, when said transformed cells are exposed to forward selection conditions the gene products of said first and second binding marker genes are deleterious or lethal to the recipient cell.
- 5
4. The vector of claim 3 in which at least one of the operons confers a genotype selected from the group consisting of galT<sup>+</sup>, tetA<sup>+</sup>, lacZ<sup>+</sup>, pheS<sup>+</sup>, argP<sup>+</sup>, thyA<sup>+</sup>, crp<sup>+</sup>, pyrF<sup>+</sup>, ptsM<sup>+</sup>, secA<sup>+</sup>, malE<sup>+</sup>, lacZ<sup>+</sup>, ompA<sup>+</sup>, btuB<sup>+</sup>, lamB<sup>+</sup>, tonA<sup>+</sup>, cir<sup>+</sup>, tsx<sup>+</sup>, aroP<sup>+</sup>, cysK<sup>+</sup>, and dctA<sup>+</sup>.
- 10
5. The vector of claim 3 wherein the binding marker genes are functionally unrelated.
6. The vector of claim 5 wherein the first and second operons confer, respectively, a pair of genotypes selected from the group consisting of:
- 15
- (a) galT<sup>+</sup> and tetA<sup>+</sup>;  
(b) argP<sup>+</sup> and pheS<sup>+</sup>;  
(c) lacZ<sup>+</sup> and tetA<sup>+</sup>;  
(d) dctA<sup>+</sup> and cysK<sup>+</sup>;  
(e) crp<sup>+</sup> and thyA<sup>+</sup>;  
(f) lamB<sup>+</sup> and thyA<sup>+</sup>;  
(g) SecA<sup>+</sup> / malE<sup>+</sup> / lacZ<sup>+</sup> and pyrF<sup>+</sup>;  
(h) tsx<sup>+</sup> and cysK<sup>+</sup>;  
(i) dctA<sup>+</sup> and thyA<sup>+</sup>;  
(j) galT<sup>+</sup> and pheS<sup>+</sup>;  
(k) tetA<sup>+</sup> and thyA<sup>+</sup>;  
(l) ptsM<sup>+</sup> and thyA<sup>+</sup>;  
(m) ompA<sup>+</sup> and pyrF<sup>+</sup>;  
(n) btuB<sup>+</sup> and pyrF<sup>+</sup>;  
(o) tonA<sup>+</sup> and galT<sup>+</sup>;  
(p) cir<sup>+</sup> and cysK<sup>+</sup>; and  
(q) aroP<sup>+</sup> and lacZ<sup>+</sup>.
- 20
- 25
- 30
- 35
7. The vector of any of claims 3-6 wherein the promoters of said first and second operons are different.
8. The vector of claim 7 wherein the degree of homology between the first and second promoters is less than 50% in the region between the -10 region of the promoter and the base at which transcription is initiated.
- 40
9. The vector of any of claims 3-8 wherein at least one of said operons comprises a plurality of copies of the target DNA sequences, wherein each copy is positioned so that the target DNA sequence interferes substantially with expression if and only if a protein expressed by the recipient cell binds to the target DNA sequence.
10. The vector of any of claims 3-9 further comprising a plurality of genetic elements essential to the maintenance of the vector or the survival of the transformed cells under conditions that select for presence of said vector, said operons and said genetic elements being positioned on said vector so no single deletion event can render nonfunctional more than one of said operons without also rendering nonfunctional one of said essential genetic elements.
- 45
11. The vector of claim 10 wherein at least one of said genetic elements comprises a selectable beneficial or essential gene, and a control promoter operably linked to said beneficial or conditionally essential gene, but where no instance of said target DNA sequence is associated with said genetic element.
- 50
12. The vector of claim 11 wherein the control promoter is essentially identical to the promoter of one of said selectable binding marker operons, so that proteins binding to the latter promoter will also bind to the control promoter and thereby inhibit expression of said beneficial or essential gene.
- 55
13. The vector of any of claims 3-12 wherein under reverse selection conditions the gene products of said binding marker genes are beneficial or conditionally essential to the transformed cells.

14. The vector of claim 13 wherein each of the first and second operons confers a phenotype selected independently but non-identically from the group consisting of: *galT.K<sup>+</sup>*, *tetA<sup>+</sup>*, *lacZ<sup>+</sup>*, *pheS<sup>+</sup>*, *argP<sup>+</sup>*, *thyA<sup>+</sup>*, *cro<sup>+</sup>*, *pyrF<sup>+</sup>*, *ptsM<sup>+</sup>*, *secA<sup>+</sup>*, *malE<sup>+</sup>*, *lacZ<sup>+</sup>*, *ompA<sup>+</sup>*, *btuB<sup>+</sup>*, *lamB<sup>+</sup>*, *tonA<sup>+</sup>*, *cir<sup>+</sup>*, *tsx<sup>+</sup>*, *aroP<sup>+</sup>*, *cysK<sup>+</sup>*, and *dctA<sup>+</sup>*.

15. The vector of any of claims 3-14 further comprising a nondeleterious cloning site so positioned that insertion of a foreign gene at such site does not inactivate said first or second operons or any genetic element of said vector required for its maintenance within the transformed cell.

16. The selection vector of any of claims 3-15 in which a target sequence associated with at least one of said operons is positioned within the RNA-polymerase binding site of the promoter of the operon.

17. The selection vector of any of claims 3-16 in which a target sequence associated with at least one of said operons is positioned upstream of the -35 region of the promoter of the operon.

18. The selection vector of any of claims 3-17 in which a target sequence associated with at least one of said operons is positioned downstream of the -10 region of the promoter of the operon.

19. The selection vector of any of claims 3-18 in which a target sequence is positioned so that the most 5' base of the target sequence is transcribed into the +1 base or the +2 base of the mRNA transcribed under the direction of the promoter of the operon.

20. The selection vector of claim 10 in which one of said genetic elements is the origin of replication of said vector.

21. The vector of any of claims 3-20 further comprising a gene (*pdbp*) coding on expression for a potential DNA-binding protein or polypeptide, said gene comprising:

- a) a coding region that codes on expression for a polypeptide, each domain of said polypeptide having at least 50% sequence identity to a known DNA-binding domain, and
- b) a promoter operably linked to said coding region for controlling its expression.

22. The vector of any of claims 3-20 wherein said first or second promoter is an inducible or repressible promoter.

23. The vector of any of claims 3-22 wherein the target DNA sequence comprises 10-25 base pairs.

24. The vector of any of claims 3-22 wherein no copy of the target DNA sequence occurs naturally in said first promoter, said second promoter, noncoding regions of said first binding marker gene or noncoding regions of said second binding marker gene.

25. The vector of claim 1 wherein the selective advantage is the ability to better utilize a particular nutrient or reduced dependency on a particular nutrient.

26. The vector of claim 1 wherein the selective advantage is resistance to a substance otherwise toxic to the recipient cell.

27. The vector of claim 26 wherein the selective advantage is resistance to an antibiotic.

28. A population of vectors according to claim 21 randomly mutated to potentially encode any of 2 to 20 predetermined amino acids, at predetermined codons within the *pdbp* gene so that said vectors collectively can express a plurality of different but sequence-related potential DNA-binding proteins.

29. The population of claim 28 wherein the level of random mutation is such that from  $10^4$  to  $10^9$  different potential DNA-binding proteins can be expressed.

30. A cell culture comprising a plurality of cells, each cell bearing a selection vector according to any of claims 1-27, said cell bearing a gene coding on expression for a potential DNA-binding protein, where said cells collectively can express a plurality of different but sequence-related potential DNA-binding proteins.

31. A method of obtaining a gene coding on expression for a novel DNA-binding protein or polypeptide that specifically

binds a predetermined DNA target sequence in double stranded DNA, comprising:

- (a) providing a cell culture according to claim 30;
  - (b) causing the cells of such culture to express said potential DNA-binding proteins or polypeptides;
  - (c) exposing the cells to forward selection conditions to select for cells which express a protein or polypeptide which specifically binds to a said target DNA sequence; and
  - (d) recovering the selected cells bearing a gene coding on expression for such protein or polypeptide.
32. A method of producing a DNA-binding protein or polypeptide which specifically binds a predetermined double stranded DNA target, which comprises obtaining a gene by the method of claim 31 which codes on expression for such protein or polypeptide, expressing the gene in a suitable host cell, and recovering said protein or polypeptide.
  33. A method of obtaining a protein or polypeptide which may be used to specifically repress a coding or regulatory element of interest which comprises obtaining a gene by the method of claim 31, determining the sequence of at least the DNA-binding domain of the protein or polypeptide, and producing at least the DNA-binding domain of the protein or polypeptide.
  34. The method of claim 31 wherein a gene encoding a known DNA binding protein picked from the group consisting of Cro from phage  $\lambda$ , cI repressor from phage  $\lambda$ , Cro from phage 434, cI repressor from phage 434, P22 repressor *E. coli* tryptophan repressor, *E. coli* CAP, P22 Arc, P22 Mnt, *E. coli* lactose repressor, MAT-a1-alpha2 from yeast, Polyoma Large T antigen, SV40 Large T antigen, Adenovirus E1A, and TFIIIA from *Xenopus laevis* is randomly mutated to potentially encode any of 2 to 20 predetermined amino acids, at predetermined codons, to obtain genes coding on expression for a plurality of potential target DNA-binding proteins.
  35. The method of claim 31 in which the DNA binding protein comprises a plurality of zinc finger DNA-binding domains.
  36. The method of claim 31 in which the cells prior to transformation are of a  $\text{Gal}^E$ ,  $\text{Gal}^T$ ,  $\text{Gal}^K$ ,  $\text{Tet}^S$  phenotype, the binding marker genes are the *tet* and *galT.K* genes, and the forward selection condition is cultivation of the cells in a medium containing galactose, or fusaric acid, or both, or substances metabolized into or catalyzing the production of galactose or fusaric acid.
  37. The method of claim 31 wherein said *pdbp* gene is randomly mutated to potentially encode any of 2 to 20 predetermined amino acids, at predetermined codons, such that at least one randomly mutated codon of said gene genetically encodes all twenty amino acids and yields at that codon a ratio of abundance of least favored amino acid to most favored amino acid which is greater than that obtained with an NNN codon, N denoting an equimolar mixture of G, A, T and C.
  38. The method of claim 37 wherein for said codon, the frequencies with which acidic and basic amino acids are encoded are equal.
  39. The method of claim 37 wherein said randomly mutated codon has substantially the following base proportions:

|         | I    | C    | A    | G    |
|---------|------|------|------|------|
| base #1 | 0.26 | 0.18 | 0.26 | 0.30 |
| base #2 | 0.22 | 0.16 | 0.40 | 0.22 |
| base #3 | 0.5  | 0.0  | 0.0  | 0.5  |

#### Patentansprüche

1. Selektionsvektor zur Selektion von Empfängerzellen, die mit einem solchen Vektor transformiert sind, welche ein Protein oder Polypeptid exprimieren, das spezifisch an eine vorbestimmte Ziel-DNA-Sequenz bindet, die von dem Vektor getragen wird, umfassend ein erstes und zweites Operon, wobei jedes mindestens ein exprimierbares Gen umfaßt, wobei die Gene des ersten und zweiten Operons unterschiedlich sind, wobei eine Kopie der Ziel-DNA-Sequenz in jedem Operon enthalten ist und darin so positioniert ist, daß die Empfängerzellen einen Selektionsvor-

teil haben, der nicht die Resistenz gegenüber lytischem Wachstum von Phagen ist, wenn Sie ein Protein oder Polypeptid exprimieren, das an die Kopien der Ziel-DNA-Sequenz bindet.

2. Vektor nach Anspruch 1, wobei mindestens ein Operon ein selektierbares nützliches Gen, einen okkludierbaren Promotor, der funktionell mit dem nützlichen Gen verknüpft ist und seine Transkription lenkt, einen okkludierenden Promotor, der die Transkription von dem okkludierbaren Promotor des nützlichen Gens okkludiert, und eine Kopie der Ziel-DNA-Sequenz umfaßt, die so positioniert ist, daß die Bindung des Proteins oder Polypeptids an die Kopie den okkludierenden Promotor reprimiert und so die Transkription des nützlichen Gens erleichtert.

3. Selektionsvektor nach Anspruch 1 oder 2, umfassend:

a) ein erstes Operon, wobei das Operon umfaßt:

- i) ein erstes Bindungsmarkergen oder -gene,
- ii) einen ersten Promotor, der die Expression des (der) Bindungsmarkergens oder -gene steuert, und
- iii) eine erste Kopie der Ziel-DNA-Sequenz, wobei die Ziel-DNA-Sequenz die Expression des (der) ersten Gens (Gene) erheblich beeinflusst, wenn ein Protein, das von der Empfängerzelle exprimiert wird, an die Ziel-DNA-Sequenz bindet,

b) ein zweites Operon, wobei das Operon umfaßt:

- i) ein zweites Bindungsmarkergen oder -gene,
  - ii) einen zweiten Promotor, der die Expression des (der) zweiten Bindungsmarkergens oder -gene steuert, und
  - iii) eine zweite Kopie der Ziel-DNA-Sequenz, wobei die Ziel-DNA-Sequenz die Expression des (der) zweiten Gens (Gene) erheblich beeinflusst, wenn ein Protein, das von der Empfängerzelle exprimiert wird, an die Ziel-DNA-Sequenz bindet,
- wobei die Bindungsmarkergene des ersten und zweiten Operons unterschiedlich sind und wobei, wenn die transformierten Zellen Vorwärts-Selektionsbedingungen ausgesetzt werden, die Genprodukte der ersten und zweiten Bindungsmarkergene schädlich oder tödlich für die Empfängerzelle sind.

4. Vektor nach Anspruch 3, in dem mindestens eines der Operons einen Genotyp verleiht, ausgewählt aus der Gruppe galT.K<sup>+</sup>, tetA<sup>+</sup>, lacZ<sup>+</sup>, pheS<sup>+</sup>, argP<sup>+</sup>, thyA<sup>+</sup>, crp<sup>+</sup>, pyrF<sup>+</sup>, ptsM<sup>+</sup>, secA<sup>+</sup>, malE<sup>+</sup>, lacZ<sup>+</sup>, ompA<sup>+</sup>, btuB<sup>+</sup>, lamB<sup>+</sup>, tonA<sup>+</sup>, cir<sup>+</sup>, tsx<sup>+</sup>, aroP<sup>+</sup>, cysK<sup>+</sup>, und dctA<sup>+</sup>.

5. Vektor nach Anspruch 3, wobei die Bindungsmarkergene in funktioneller Hinsicht nicht verwandt sind.

6. Vektor nach Anspruch 5, wobei das erste bzw. zweite Operon ein Genotyp-Paar verleihen, ausgewählt aus der Gruppe

- (a) galT.K<sup>+</sup> und tetA<sup>+</sup>;
- (b) argP<sup>+</sup> und pheS<sup>+</sup>;
- (c) lacZ<sup>+</sup> und tetA<sup>+</sup>;
- (d) dctA<sup>+</sup> und cysK<sup>+</sup>;
- (e) crp<sup>+</sup> und thyA<sup>+</sup>;
- (f) lamB<sup>+</sup> und thyA<sup>+</sup>;
- (g) SecA<sup>+</sup> / malE<sup>+</sup> / lacZ<sup>+</sup> und pyrF<sup>+</sup>;
- (h) tsx<sup>+</sup> und cysK<sup>+</sup>;
- (i) dctA<sup>+</sup> und thyA<sup>+</sup>;
- (j) galT.K<sup>+</sup> und pheS<sup>+</sup>;
- (k) tetA<sup>+</sup> und thyA<sup>+</sup>;
- (l) ptsM<sup>+</sup> und thyA<sup>+</sup>;
- (m) ompA<sup>+</sup> und pyrF<sup>+</sup>;
- (n) btuB<sup>+</sup> und pyrF<sup>+</sup>;
- (o) tonA<sup>+</sup> und galT.K<sup>+</sup>;
- (p) cir<sup>+</sup> und cysK<sup>+</sup>; und
- (q) aroP<sup>+</sup> und lacZ<sup>+</sup>.



7. Vektor nach einem der Ansprüche 3 bis 6, wobei die Promotoren der ersten und zweiten Operons unterschiedlich sind.
8. Vektor nach Anspruch 7, wobei der Grad der Homologie zwischen den ersten und zweiten Promotoren im Bereich zwischen dem -10-Bereich des Promotors und der Base, bei der die Transkription beginnt, geringer als 50 % ist.
9. Vektor nach einem der Ansprüche 3 bis 8, wobei mindestens eines der Operons eine Vielzahl von Kopien der Ziel-DNA-Sequenzen umfaßt, wobei jede Kopie so positioniert ist, daß die Ziel-DNA-Sequenz die Expression erheblich beeinflusst, wenn und nur wenn ein Protein, das von der Empfängerzelle exprimiert wird, an die Ziel-DNA-Sequenz bindet.
10. Vektor nach einem der Ansprüche 3 bis 9, ferner umfassend eine Vielzahl von genetischen Elementen, die wesentlich sind für die Erhaltung des Vektors oder das Überleben der transformierten Zellen unter Bedingungen, die auf das Vorhandensein des Vektors selektieren, wobei die Operons und die genetischen Elemente so auf dem Vektor positioniert sind, daß kein Einzeldelotionsereignis eine Nichtfunktionalität von mehr als einem der Operons verursachen kann, ohne auch eines der wesentlichen genetischen Elemente nichtfunktionell zu machen.
11. Vektor nach Anspruch 10, wobei mindestens eines der genetischen Elemente ein selektierbar nützliches oder essentielles Gen umfaßt, und einen Kontrollpromotor, der funktionell mit dem nützlichen oder gegebenenfalls essentiellen Gen verknüpft ist, wobei aber keinesfalls die Ziel-DNA-Sequenz mit dem genetischen Element assoziiert ist.
12. Vektor nach Anspruch 11, wobei der Kontrollpromotor im Wesentlichen identisch mit dem Promotor von einem der selektierbaren Bindungsmarkeroperons ist, so daß Proteine, die an den letzteren Promotor binden, auch an den Kontrollpromotor binden und dabei die Expression des nützlichen oder essentiellen Gens hemmen.
13. Vektor nach einem der Ansprüche 3 bis 12, wobei unter reversen Selektionsbedingungen die Genprodukte der Bindungsmarkergene nützlich oder bedingt essentiell für die transformierten Zellen sind.
14. Vektor nach Anspruch 13, wobei jedes der ersten und zweiten Operons einen Phänotypen verleiht, unabhängig aber nicht identisch ausgewählt aus der Gruppe bestehend aus: *galT*<sup>+</sup>, *tetA*<sup>+</sup>, *lacZ*<sup>+</sup>, *pheS*<sup>+</sup>, *argP*<sup>+</sup>, *thyA*<sup>+</sup>, *crp*<sup>+</sup>, *pyrF*<sup>+</sup>, *ptsM*<sup>+</sup>, *secA*<sup>+</sup>, *malE*<sup>+</sup>, *lacZ*<sup>+</sup>, *ompA*<sup>+</sup>, *btuB*<sup>+</sup>, *lamB*<sup>+</sup>, *tonA*<sup>+</sup>, *qir*<sup>+</sup>, *tsx*<sup>+</sup>, *aroP*<sup>+</sup>, *cysK*<sup>+</sup>, und *dctA*<sup>+</sup>.
15. Vektor nach einem der Ansprüche 3 bis 14, ferner umfassend eine nicht-schädliche Clonierungsstelle, die so positioniert ist, daß die Insertion eines Fremdgens an der Stelle nicht die ersten oder zweiten Operons oder irgendein genetisches Element des Vektors, das für die Erhaltung innerhalb der transformierten Zelle erforderlich ist, inaktiviert.
16. Selektionsvektor nach einem der Ansprüche 3 bis 15, wobei eine Zielsequenz, die mit mindestens einem der Operons assoziiert ist, innerhalb der RNA-Polymerase-Bindungsstelle des Promotors des Operons positioniert ist.
17. Selektionsvektor nach einem der Ansprüche 3 bis 16, wobei eine Zielsequenz, die mit mindestens einem der Operons assoziiert ist, stromaufwärts vom -35-Bereich des Promotors des Operons positioniert ist.
18. Selektionsvektor nach einem der Ansprüche 3 bis 17, wobei eine Zielsequenz, die mit mindestens einem der Operons assoziiert ist, stromabwärts vom -10-Bereich des Promotors des Operons positioniert ist.
19. Selektionsvektor nach einem der Ansprüche 3 bis 18, wobei eine Zielsequenz so positioniert ist, daß die 5'-nächstgelegene Base der Zielsequenz in die +1-Base oder +2-Base der mRNA transkribiert wird, die unter der Steuerung des Promotors des Operons transkribiert wird.
20. Selektionsvektor nach Anspruch 10, wobei eines der genetischen Elemente der Replikationsursprung des Vektors ist.
21. Vektor nach einem der Ansprüche 3 bis 20, ferner umfassend ein Gen (*pdbp*), das bei Expression ein potentielles DNA-Bindungsprotein oder -polypeptid codiert, wobei das Gen umfaßt:

a) einen Codierungsbereich, der bei Expression ein Polypeptid codiert, wobei jede Domäne des Polypeptids

mindestens 50 % Sequenzidentität zu einer bekannten DNA-Bindungsdomäne aufweist, und  
b) einen Promotor, der funktionell verknüpft ist mit dem Codierungsbereich zur Kontrolle seiner Expression.

22. Vektor nach einem der Ansprüche 3 bis 20, wobei der erste oder zweite Promotor ein induzierbarer oder reprimierbarer Promotor ist.
23. Vektor nach einem der Ansprüche 3 bis 22, wobei die Ziel-DNA-Sequenz 10 bis 25 Basenpaare umfaßt.
24. Vektor nach einem der Ansprüche 3 bis 22, wobei keine Kopie der Ziel-DNA-Sequenz natürlicherweise in dem ersten Promotor, dem zweiten Promotor, den nicht-codierenden Bereichen des ersten Bindungsmerkergens oder den nichtcodierenden Bereichen des zweiten Bindungsmerkergens vorkommt.
25. Vektor nach Anspruch 1, wobei der Selektionsvorteil die Fähigkeit ist, ein bestimmtes Nahrungsmittel besser zu verwerten, oder eine geringere Abhängigkeit von einem bestimmten Nahrungsmittel.
26. Vektor nach Anspruch 1, wobei der Selektionsvorteil die Resistenz gegenüber einer Substanz ist, die sonst toxisch für die Empfängerzelle ist.
27. Vektor nach Anspruch 26, wobei der Selektionsvorteil die Resistenz gegenüber einem Antibiotikum ist.
28. Population von Vektoren nach Anspruch 21, die zufällig mutiert sind, so daß sie möglicherweise eine der 2 bis 20 vorbestimmten Aminosäuren codieren, an vorbestimmten Codons innerhalb des *pdbp*-Gens, so daß die Vektoren vereint eine Vielzahl von unterschiedlichen aber sequenzverwandten möglichen DNA-Bindungsproteinen exprimieren können.
29. Population nach Anspruch 28, wobei das Ausmaß der Zufallsmutation so ist, daß  $10^4$  bis  $10^9$  unterschiedliche potentielle DNA-Bindungsproteinen exprimiert werden können.
30. Zellkultur, umfassend eine Vielzahl von Zellen, wobei jede Zelle einen Selektionsvektor nach einem der Ansprüche 1 bis 27 trägt, wobei die Zelle ein Gen trägt, das bei Expression ein potentielles DNA-Bindungsprotein codiert, wobei die Zellen vereint eine Vielzahl von unterschiedlichen aber sequenzverwandten potentiellen DNA-Bindungsproteinen exprimieren können.
31. Verfahren zum Erhalt eines Gens, das bei Expression ein neues DNA-Bindungsprotein oder -polypeptid codiert, das spezifisch an eine vorbestimmte DNA-Zielsequenz in doppelsträngiger DNA bindet, umfassend:
  - a) Bereitstellung einer Zellkultur nach Anspruch 30,
  - b) Veranlassung der Zellen einer solchen Kultur, die potentiellen DNA-Bindungsproteine oder -polypeptide zu exprimieren,
  - c) Unterwerfen der Zellen unter Vorwärts-Selektionsbedingungen, um Zellen auszuwählen, die ein Protein oder Polypeptid exprimieren, das spezifisch an die Ziel-DNA-Sequenz bindet, und
  - d) Gewinnung der ausgewählten Zellen, die ein Gen tragen, das bei Expression ein solches Protein oder Polypeptid codiert.
32. Verfahren zur Produktion eines DNA-Bindungsproteins oder -polypeptids, das spezifisch an ein vorbestimmtes doppelsträngiges DNA-Ziel bindet, umfassend das Erhalten eines Gens durch das Verfahren nach Anspruch 31, welches bei Expression ein solches Protein oder Polypeptid codiert, Expression des Gens in einer geeigneten Wirtszelle und Gewinnung des Proteins oder Polypeptids.
33. Verfahren zum Erhalten eines Proteins oder Polypeptids, das verwendet werden kann, um spezifisch ein codierendes oder regulatorisches Element von Interesse zu reprimieren, umfassend das Erhalten eines Gens gemäß dem Verfahren von Anspruch 31, Bestimmung der Sequenz von mindestens der DNA-Bindungsdomäne des Proteins oder Polypeptids und Produktion von mindestens einer DNA-Bindungsdomäne des Proteins oder Polypeptids.
34. Verfahren nach Anspruch 31, wobei ein Gen, das ein bekanntes DNA-Bindungsprotein codiert, ausgewählt aus der Gruppe bestehend aus Cro vom Phagen  $\lambda$ ,  $\text{cl}$ -Repressor vom Phagen  $\lambda$ , Cro vom Phagen 434,  $\text{cl}$ -Repressor vom Phagen 434, P22-Repressor, *E. coli*-Tryptophan-Repressor, *E. coli*-CAP, P22 Arc, P22 Mnt, *E. coli*-Lactose-Repressor, MAT-a1-alpha2 aus Hefe, Polyoma Large T-Antigen, SV40 Large T-Antigen, Adenovirus E1A, und TFI-

IIA von *Xenopus laevis* zufällig mutiert ist, um möglicherweise eine der 2 bis 20 vorbestimmten Aminosäuren zu codieren, an vorbestimmten Codons, so daß Gene erhalten werden, die bei Expression eine Vielzahl von möglichen Ziel-DNA-Bindungsproteinen codieren.

- 5 35. Verfahren nach Anspruch 31, wobei das DNA-Bindungsprotein eine Vielzahl von Zink-Finger-DNA-Bindungsdomänen umfaßt.
36. Verfahren nach Anspruch 31, wobei die Zellen vor der Transformation vom GalE<sup>-</sup>, GalT<sup>-</sup>, GalK<sup>-</sup>, Tet<sup>S</sup>-Phenotyp sind, die Bindungsmarkergene die *tet*- und *galT,K*-Gene sind und die Vorwärts-Selektionsbedingung die Züchtung  
10 der Zellen in einem Medium ist, das Galactose oder Fusarinsäure oder beides enthält oder Substanzen, die in Galactose oder Fusarinsäure verstoffwechselt werden oder deren Produktion katalysieren.
37. Verfahren nach Anspruch 31, wobei das *pdbp*-Gen zufällig mutiert ist, so daß es eine der 2 bis 20 vorbestimmten Aminosäuren codiert, an vorbestimmten Codons, so daß mindestens ein zufällig mutiertes Codon des Gens genetisch alle 20 Aminosäuren codiert und an diesem Codon einen Überschußanteil der am wenigsten bevorzugten Aminosäure bis zur am meisten bevorzugten Aminosäure aufweist, der größer ist als der, der mit einem NNN-Codon erhalten wird, wobei N ein equimolares Gemisch von G, A, T und C bedeutet.
- 15 38. Verfahren nach Anspruch 37, wobei für das Codon die Häufigkeiten, mit denen saure und basische Aminosäuren codiert werden, gleich sind.
39. Verfahren nach Anspruch 37, wobei das zufällig mutierte Codon im Wesentlichen die folgenden Basenanteile aufweist:

|            | I    | C    | A    | G    |
|------------|------|------|------|------|
| Base Nr. 1 | 0.26 | 0.18 | 0.26 | 0.30 |
| Base Nr. 2 | 0.22 | 0.16 | 0.40 | 0.22 |
| Base Nr. 3 | 0.5  | 0.0  | 0.0  | 0.5  |

#### Revendications

- 35 1. Vecteur de sélection pour sélectionner des cellules réceptrices transformées par un tel vecteur qui expriment une protéine ou un polypeptide qui se lie spécifiquement à une borne de séquence d'ADN cible prédéterminée par ledit vecteur, qui comprend un premier et un deuxième opéron, comprenant chacun au moins un gène exprimable, les gènes desdits premier et deuxième opérons étant différents, une copie de la séquence d'ADN cible étant incluse  
40 dans chaque opéron et positionnée dedans de sorte que les cellules réceptrices jouissent d'un avantage sélectif, autre qu'une résistance à la croissance lytique de phage, si elles peuvent exprimer une protéine ou un polypeptide qui se lie auxdites copies de la séquence d'ADN cible.
- 45 2. Vecteur selon la revendication 1, dans lequel au moins un opéron comprend un gène bénéfique sélectionnable, un promoteur pouvant être bloqué lié de manière utilisable audit gène bénéfique et dirigeant sa transcription, un promoteur de blocage bloquant la transcription à partir dudit promoteur pouvant être bloqué dudit gène bénéfique, et une copie de la séquence d'ADN cible positionnée de sorte que la liaison de ladite protéine ou du polypeptide à ladite copie réprime ledit promoteur de blocage et facilite ainsi la transcription dudit gène bénéfique.
- 50 3. Vecteur de sélection selon la revendication 1 ou 2, qui comprend:
- a) un premier opéron, opéron qui comprend:
- 55 i) un ou des premiers gènes marqueurs de liaison;
- ii) un premier promoteur régulant l'expression dudit ou desdits gènes marqueurs de liaison; et
- iii) une première copie de la séquence d'ADN cible, où ladite séquence d'ADN cible interfère sensiblement avec l'expression du ou des premiers gènes si une protéine exprimée par la cellule réceptrice se lie à la séquence d'ADN cible;

b) un deuxième opéron, opéron qui comprend:

- i) un ou des deuxièmes gènes marqueurs de liaison;
  - ii) un deuxième promoteur régulant l'expression dudit ou desdits deuxièmes gènes marqueurs de liaison;
  - et
  - iii) une deuxième copie de la séquence d'ADN cible, où ladite séquence d'ADN cible interfère sensiblement avec l'expression du ou des deuxièmes gènes si une protéine exprimée par la cellule réceptrice se lie à la séquence d'ADN cible,
- où les gènes marqueurs de liaison desdits premier et deuxième opérons sont différents, et où, quand lesdites cellules transformées sont exposées à des conditions de sélection avancées, les produits des gènes desdits premier et deuxième gènes marqueurs de liaison sont délétères ou léthaux à la cellule réceptrice.

4. Vecteur selon la revendication 3, dans lequel au moins un des opérons confère un génotype choisi dans le groupe formé par galT.K<sup>±</sup>, tetA<sup>±</sup>, lacZ<sup>±</sup>, pheS<sup>±</sup>, argP<sup>±</sup>, thyA<sup>±</sup>, crp<sup>±</sup>, pyrF<sup>±</sup>, ptsM<sup>±</sup>, secA<sup>±</sup>/malE<sup>±</sup>/lacZ<sup>±</sup>, ompA<sup>±</sup>, btuB<sup>±</sup>, lamB<sup>±</sup>, tonA<sup>±</sup>, cir<sup>±</sup>, tsx<sup>±</sup>, aroP<sup>±</sup>, cysK<sup>±</sup> et dctA<sup>±</sup>.

5. Vecteur selon la revendication 3, dans lequel les gènes marqueurs de liaison ne sont pas fonctionnellement apparentés.

6. Vecteur selon la revendication 5, dans lequel les premier et deuxième opérons confèrent, respectivement, une paire de génotypes choisis dans le groupe formé par:

- (a) galT.K<sup>±</sup> et tetA<sup>±</sup>;
- (b) argP<sup>±</sup> et pheS<sup>±</sup>;
- (c) lacZ<sup>±</sup> et tetA<sup>±</sup>;
- (d) dctA<sup>±</sup> et cysK<sup>±</sup>;
- (e) crp<sup>±</sup> et thyA<sup>±</sup>;
- (f) lamB<sup>±</sup> et thyA<sup>±</sup>;
- (g) SecA<sup>±</sup>/malE<sup>±</sup>/lacZ<sup>±</sup> et pyrF<sup>±</sup>;
- (h) tsx<sup>±</sup> et cysK<sup>±</sup>;
- (i) dctA<sup>±</sup> et thyA<sup>±</sup>;
- (j) galT.K<sup>±</sup> et pheS<sup>±</sup>;
- (k) tetA<sup>±</sup> et thyA<sup>±</sup>;
- (l) ptsM<sup>±</sup> et thyA<sup>±</sup>;
- (m) ompA<sup>±</sup> et pyrF<sup>±</sup>;
- (n) btuB<sup>±</sup> et pyrF<sup>±</sup>;
- (o) tonA<sup>±</sup> et galT.K<sup>±</sup>;
- (p) cir<sup>±</sup> et cysK<sup>±</sup>; et
- (q) aroP<sup>±</sup> et lacZ<sup>±</sup>.

7. Vecteur selon l'une quelconque des revendications 3 à 6, dans lequel les promoteurs desdits premier et deuxième opérons sont différents.

8. Vecteur selon la revendication 7, dans lequel le degré d'homologie entre les premier et deuxième promoteurs est inférieur à 90% dans la zone entre la zone -10 du promoteur et la base au niveau de laquelle la transcription est amorcée.

9. Vecteur selon l'une quelconque des revendications 3 à 8, dans lequel au moins un desdits opérons comprend une pluralité de copies des séquences d'ADN cible, dans lequel chaque copie est positionnée de sorte que la séquence d'ADN cible interfère sensiblement avec l'expression si et seulement si une protéine exprimée par la cellule réceptrice se lie à la séquence d'ADN cible.

10. Vecteur selon l'une quelconque des revendications 3 à 9, comprenant de plus une pluralité d'éléments génétiques essentiels à la maintenance du vecteur ou à la survie des cellules transformées dans des conditions qui sélectionnent pour la présence dudit vecteur, lesdits opérons et lesdits éléments génétique étant positionnés sur ledit vecteur, ainsi aucun cas de délétion unique ne peut rendre non fonctionnel plus d'un desdits opérons sans rendre également non fonctionnel un desdits éléments génétiques essentiels.

11. Vecteur selon la revendication 10, dans lequel un desdits éléments génétiques comprend un gène sélectivement bénéfique ou essentiel, et un promoteur témoin lié de façon utilisable audit gène bénéfique ou conditionnellement essentiel, mais où aucun exemple de ladite séquence d'ADN cible n'est associé audit élément génétique.
- 5 12. Vecteur selon la revendication 11, dans lequel le promoteur témoin est essentiellement identique à le promoteur d'un desdits opérons marqueurs de liaison sélectionnables, de sorte que les protéines se liant au dernier promoteur se lieront également au promoteur témoin et inhiberont ainsi l'expression dudit gène bénéfique ou essentiel.
- 10 13. Vecteur selon la revendication selon l'une quelconque des revendications 3 à 12, dans lequel, dans des conditions de sélection inverse, les produits de gènes desdits gènes marqueurs de liaison sont bénéfiques ou conditionnellement essentiels aux cellules transformées.
- 15 14. Vecteur selon la revendication 13, dans lequel chacun des premier et deuxième opérons confère un phénotype choisi indépendamment mais pas de manière identique dans le groupe formé par: galT.K<sup>+</sup>, tetA<sup>+</sup>, lacZ<sup>+</sup>, pheS<sup>+</sup>, argP<sup>+</sup>, thyA<sup>+</sup>, crp<sup>+</sup>, pyrF<sup>+</sup>, ptsM<sup>+</sup>, secA<sup>+</sup>/malE<sup>+</sup>/lacZ<sup>+</sup>, ompA<sup>+</sup>, btuB<sup>+</sup>, lamB<sup>+</sup>, tonA<sup>+</sup>, cir<sup>+</sup>, tsx<sup>+</sup>, aroP<sup>+</sup>, cysK<sup>+</sup> et dctA<sup>+</sup>.
- 20 15. Vecteur selon l'une quelconque des revendications 3 à 14, comprenant de plus un site de clonage non délétère positionné de façon à ce que l'insertion d'un gène étranger au niveau d'un tel site n'inactive pas lesdits premier ou deuxième opérons ou tout élément génétique dudit vecteur nécessaire pour sa maintenance dans la cellule transformée.
- 25 16. Vecteur de sélection selon l'une quelconque des revendications 3 à 15, dans lequel une séquence cible associée à au moins un desdits opérons est positionnée dans le site de liaison d'ARN-polymérase du promoteur de l'opéron.
- 30 17. Vecteur de sélection selon l'une quelconque des revendications 3 à 16, dans lequel une séquence cible associée à au moins un desdits opérons est positionnée en amont de la zone -35 du promoteur de l'opéron.
- 35 18. Vecteur de sélection selon l'une quelconque des revendications 3 à 17, dans lequel une séquence cible associée à au moins un desdits opérons est positionnée en aval de la zone -10 du promoteur de l'opéron.
- 40 19. Vecteur de sélection selon l'une quelconque des revendications 3 à 18, dans lequel une séquence cible est positionnée de sorte que la plus grande partie de la base 5' de la séquence cible est transcrite dans la base +1 ou la base +2 de l'ARNm transcrit sous la direction du promoteur de l'opéron.
- 45 20. Vecteur de sélection selon la revendication 10, dans lequel un desdits éléments génétiques est l'origine de la répllication dudit vecteur.
- 50 21. Vecteur selon l'une quelconque des revendications 3 à 20, comprenant de plus un gène (pdibp) codant l'expression pour une protéine ou un polypeptide de liaison à un ADN potentiel, ledit gène comprenant:
- a) une région codante qui code l'expression pour un polypeptide, chaque domaine dudit polypeptide ayant au moins 50% d'identité de séquence à un domaine de liaison à un ADN connu; et
- b) un promoteur lié de façon utilisable à ladite région codante pour réguler son expression.
- 55 22. Vecteur selon l'une quelconque des revendications 3 à 20, dans lequel ledit premier ou deuxième promoteur est un promoteur inductible ou répressible.
23. Vecteur selon l'une quelconque des revendications 3 à 22, dans lequel la séquence d'ADN cible comprend 10 à 25 paires de base:
24. Vecteur selon l'une quelconque des revendications 3 à 22, dans lequel aucune copie de la séquence d'ADN cible n'apparaît naturellement dans ledit premier promoteur, ledit deuxième promoteur, les régions non codantes dudit premier gène marqueur de liaison ou les régions non codantes dudit deuxième gène marqueur de liaison.
25. Vecteur selon la revendication 1, dans lequel l'avantage sélectif est la capacité à mieux utiliser un nutriment particulier ou une dépendance réduite envers un nutriment particulier.

26. Vecteur selon la revendication 1, dans lequel l'avantage sélectif est une résistance à une substance autrement toxique envers la cellule réceptrice.

27. Vecteur selon la revendication 26, dans lequel l'avantage sélectif est une résistance à un antibiotique.

28. Population de vecteurs selon la revendication 21, mutés de façon aléatoire pour coder potentiellement l'un quelconque de 2 à 20 aminoacides prédéterminés, au niveau de codons prédéterminés dans le gène pdbp de sorte que lesdits vecteurs peuvent exprimer collectivement une pluralité de protéines de liaison à un ADN potentielles différentes mais apparentées en séquence.

29. Population selon la revendication 28, dans laquelle le niveau de mutation aléatoire est tel que de  $10^4$  à  $10^9$  protéines de liaison à un ADN potentielles différentes peuvent être exprimées.

30. Culture cellulaire comprenant une pluralité de cellules, chaque cellule portant un vecteur de sélection selon l'une quelconque des revendications 1 à 27, ladite cellule portant un gène codant l'expression pour une protéine de liaison à un ADN potentielle, où lesdites cellules peuvent collectivement exprimer une pluralité de protéines de liaison à un ADN potentielles différentes mais apparentées en séquence.

31. Procédé pour obtenir un gène codant l'expression pour une nouvelle protéine ou un polypeptide de liaison à un ADN qui se lie spécifiquement à une séquence cible d'ADN prédéterminée dans l'ADN double-brin, comprenant les étapes consistant:

(a) à fournir une culture cellulaire selon la revendication 30;

(b) à faire exprimer par les cellules d'une telle culture lesdites protéines ou lesdits polypeptides de liaison à un ADN potentiels;

(c) à exposer les cellules à des conditions de sélection avancée pour choisir des cellules qui expriment une protéine ou un polypeptide qui se lie spécifiquement à une desdites séquences d'ADN cibles; et

(d) à récupérer les cellules choisies portant un gène codant l'expression pour une telle protéine ou polypeptide.

32. Procédé pour produire une protéine ou un polypeptide de liaison à un ADN qui se lie spécifiquement à une cible d'ADN double-brin prédéterminée, qui consiste à obtenir un gène par le procédé selon la revendication 31, qui code l'expression pour une telle protéine ou polypeptide, à exprimer le gène dans une cellule hôte appropriée, et à récupérer ladite protéine ou ledit polypeptide.

33. Procédé pour obtenir une protéine ou un polypeptide qui peut être utilisé pour réprimer spécifiquement un élément de codage ou régulateur considéré qui consiste à obtenir un gène par le procédé de la revendication 31, à déterminer la séquence d'au moins le domaine de liaison à un ADN de la protéine ou du polypeptide, et à produire au moins le domaine de liaison à un ADN de la protéine ou du polypeptide.

34. Procédé selon la revendication 31, dans lequel un gène codant une protéine de liaison à un ADN connue prélevée dans le groupe formé par Cro provenant du phage  $\lambda$ , répresseur du *ci* provenant du phage  $\lambda$ , Cro provenant du phage 434, répresseur du *ci* provenant du phage 434, répresseur P22, répresseur du tryptophane de *E.coli*, *E.coli* CAP, P22 Arc, P22 Mnt, répresseur du lactose de *E.coli*, MAT-a1-alpha2 provenant de la levure, antigène Polyoma Large T, antigène SV40 Large T, Adénovirus E1A, et TFIIA provenant de *Xenopus laevis* est muté de façon aléatoire pour coder potentiellement l'un quelconque des 2 à 20 aminoacides prédéterminés, au niveau de codons prédéterminés, pour obtenir des gènes codant l'expression pour une pluralité de protéines de liaison à un ADN cible potentielles.

35. Procédé selon la revendication 31, dans lequel la protéine liant l'ADN comprend une pluralité de domaines de liaison à un ADN type doigt à zinc.

36. Procédé selon la revendication 31, dans lequel les cellules avant la transformation sont d'un phénotype Gal<sup>E</sup>, Gal<sup>T</sup>, Gal<sup>K</sup>, Tet<sup>R</sup>, les gènes marqueurs de liaison sont les gènes *tet* et *galT.K*, et la condition de sélection avancée est la culture des cellules dans un milieu contenant du galactose, ou de l'acide fusarique, ou les deux, ou des substances métabolisées à l'intérieur ou en catalysant la production de galactose ou d'acide fusarique.

37. Procédé selon la revendication 31, dans lequel ledit gène pdbp est muté de façon aléatoire pour coder potentiellement l'un quelconque des 2 à 20 aminoacides prédéterminés, au niveau de codons prédéterminés, de sorte qu'au

moins un codon muté de façon aléatoire dudit gène code génétiquement l'ensemble des vingt aminoacides et produit pour ce codon un rapport d'abondance de l'acide le moins favorisé à l'acide le plus favorisé qui est supérieur à celui obtenu avec un codon NNN, N représentant un mélange équimolaire de G, A, T et C.

- 5 38. Procédé selon la revendication 37, dans lequel pour ledit codon, les fréquences, avec lesquelles les aminoacides acides et basiques sont codés, sont égales.
39. Procédé selon la revendication 37, dans lequel ledit codon muté de façon aléatoire possède sensiblement les proportions suivantes de bases:

10

|         | T    | C    | A    | G    |
|---------|------|------|------|------|
| Base #1 | 0,26 | 0,18 | 0,26 | 0,30 |
| Base #2 | 0,22 | 0,16 | 0,40 | 0,22 |
| Base #3 | 0,5  | 0,0  | 0,0  | 0,5  |

15

20

25

30

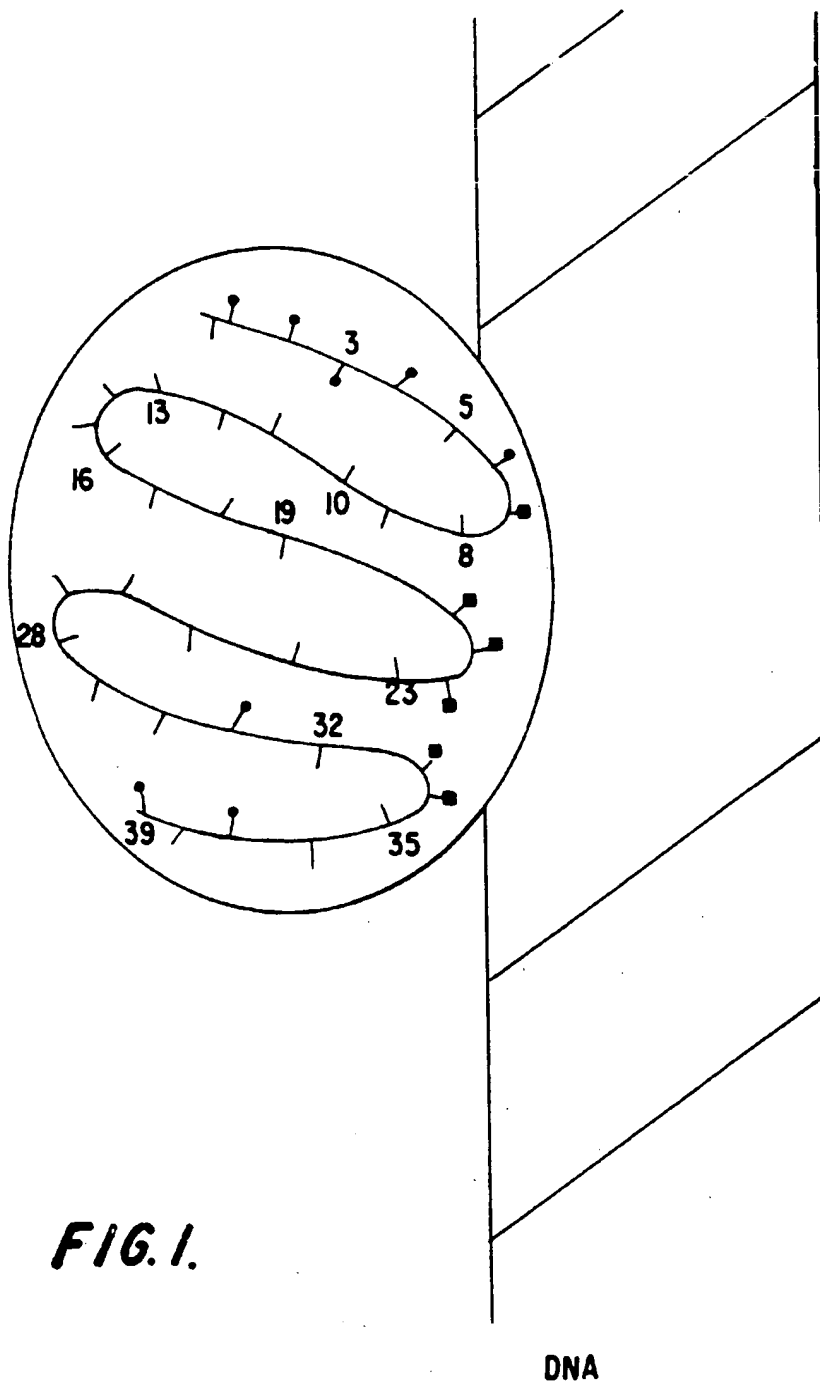
35

40

45

50

55





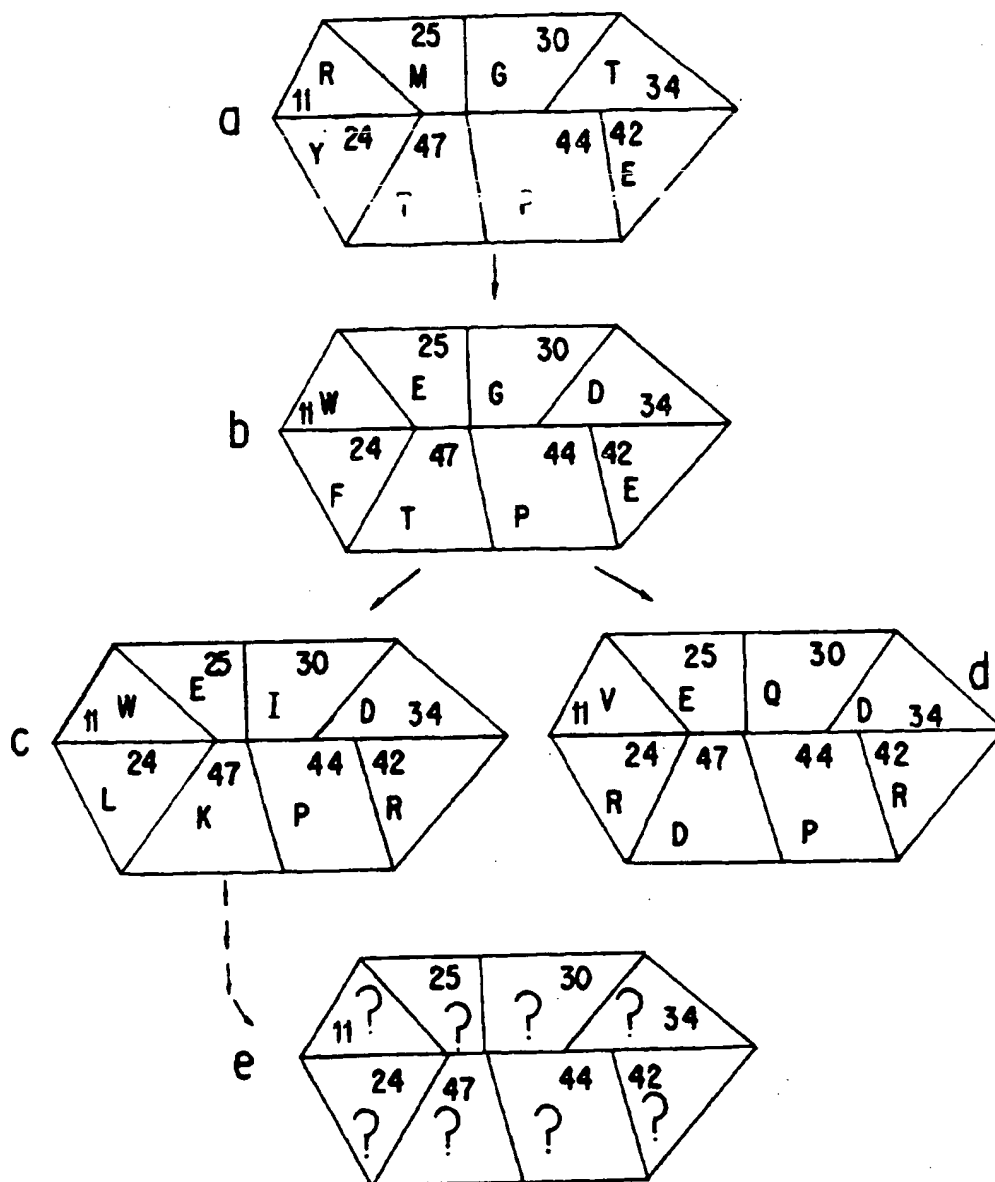
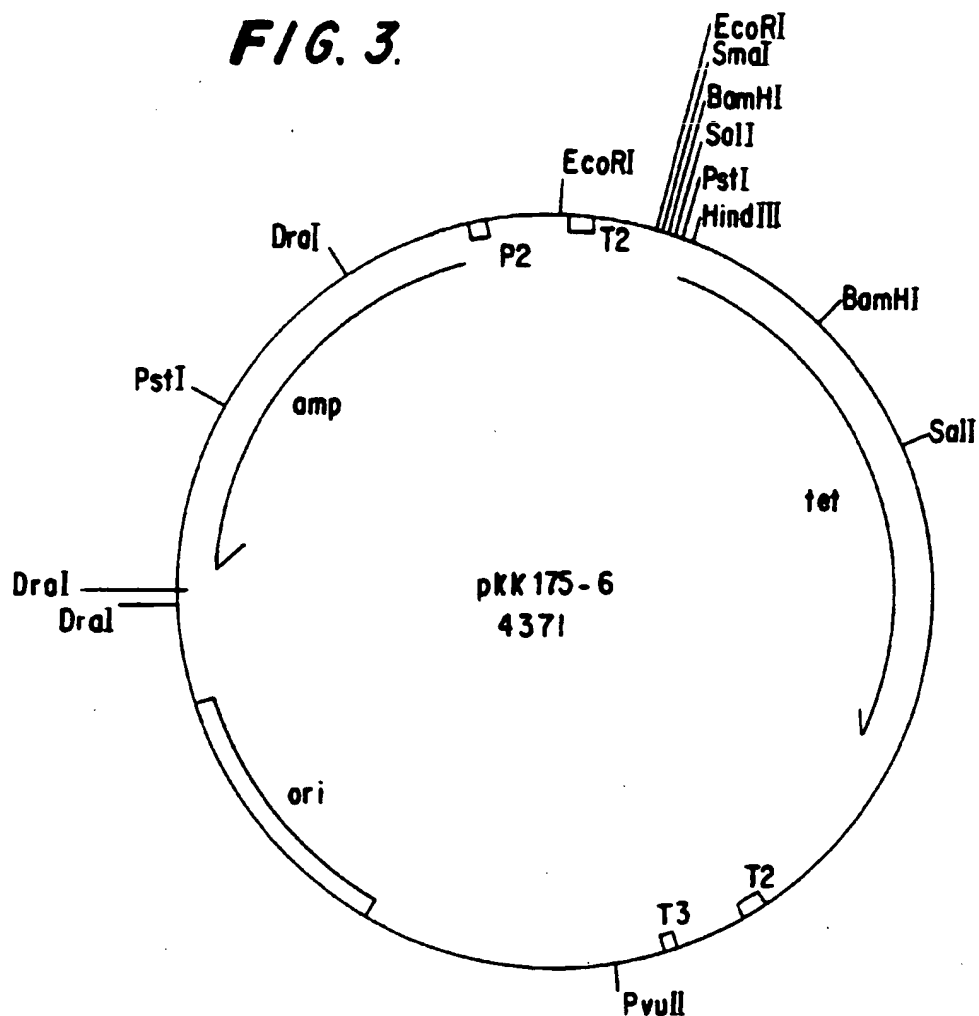


FIG. 2.

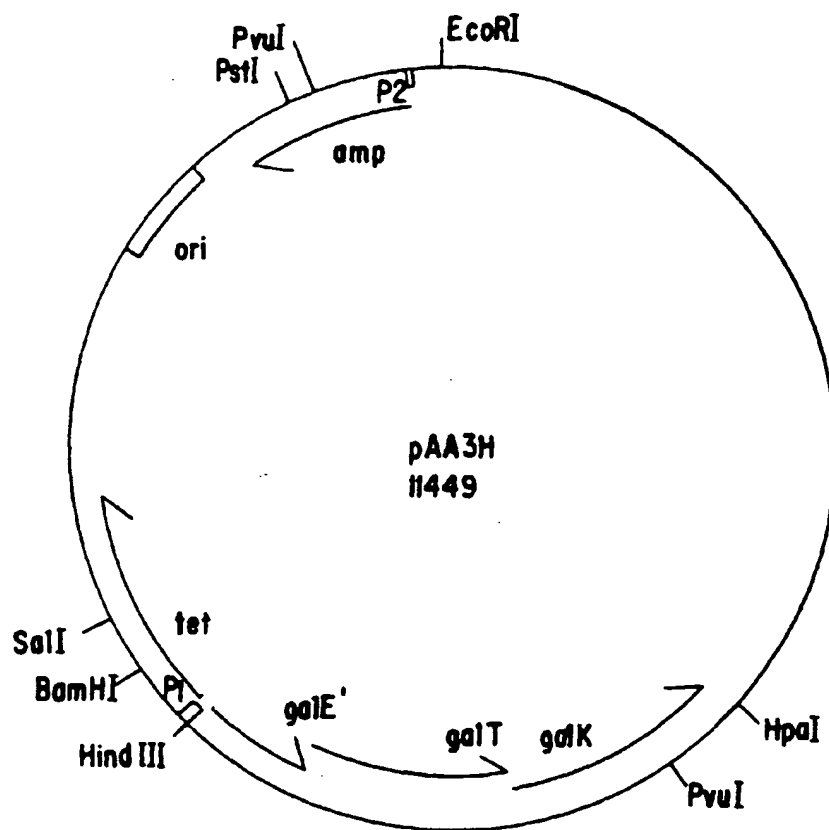
**FIG. 3.**



DELTA 4 CELLS CONTAINING pKK175-6 ARE  
Amp<sup>R</sup>, Tet<sup>S</sup>, Fus<sup>R</sup>, AND Gal<sup>R</sup>.

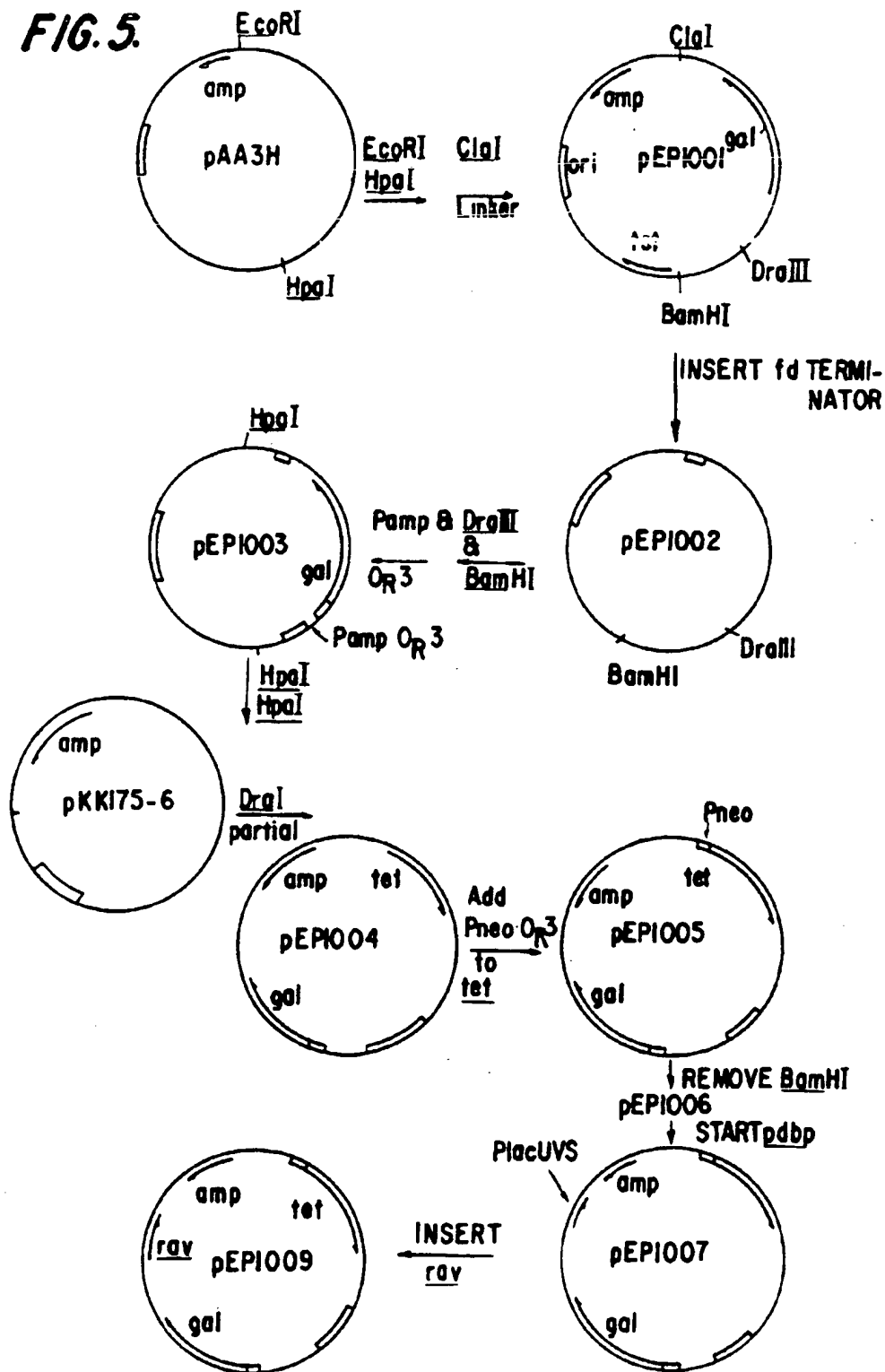
T2 IS rrnBt1; T3 IS rrnBt2; P2 IS Pamp

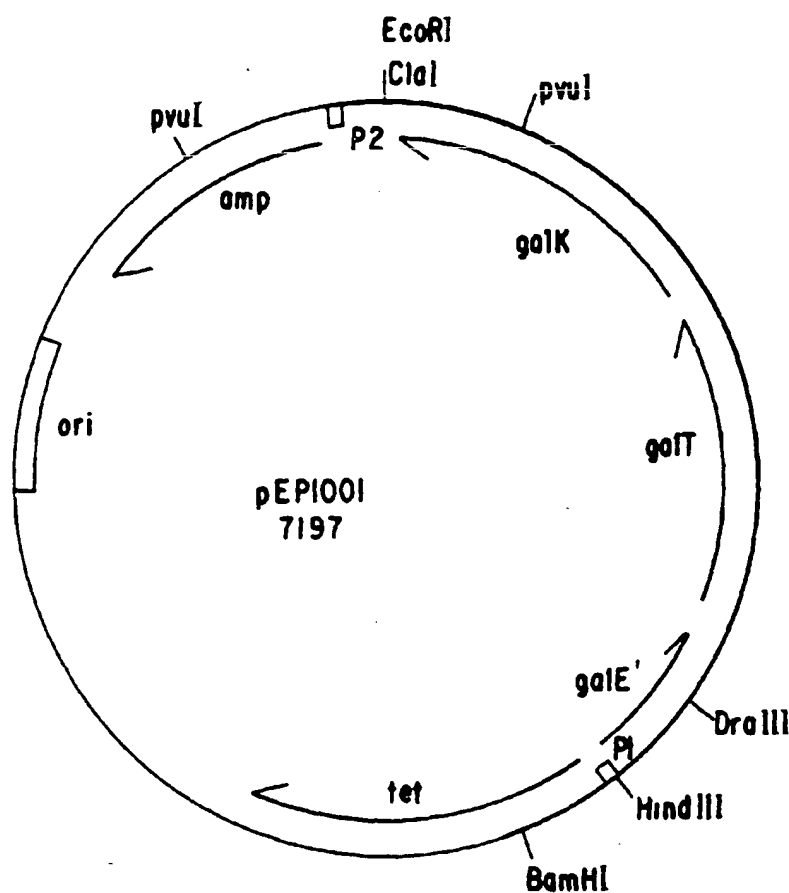
**FIG. 4.**



**DELTA4 CELLS CONTAINING pAA3H are  
Amp<sup>R</sup>, Tet<sup>S</sup>, Fus<sup>R</sup>, AND Gal<sup>S</sup>**

**T2 IS rrnB1; T3 IS rrnB2;  
PI IS pBR322 PI PROMOTER;  
P2 IS Pamp**

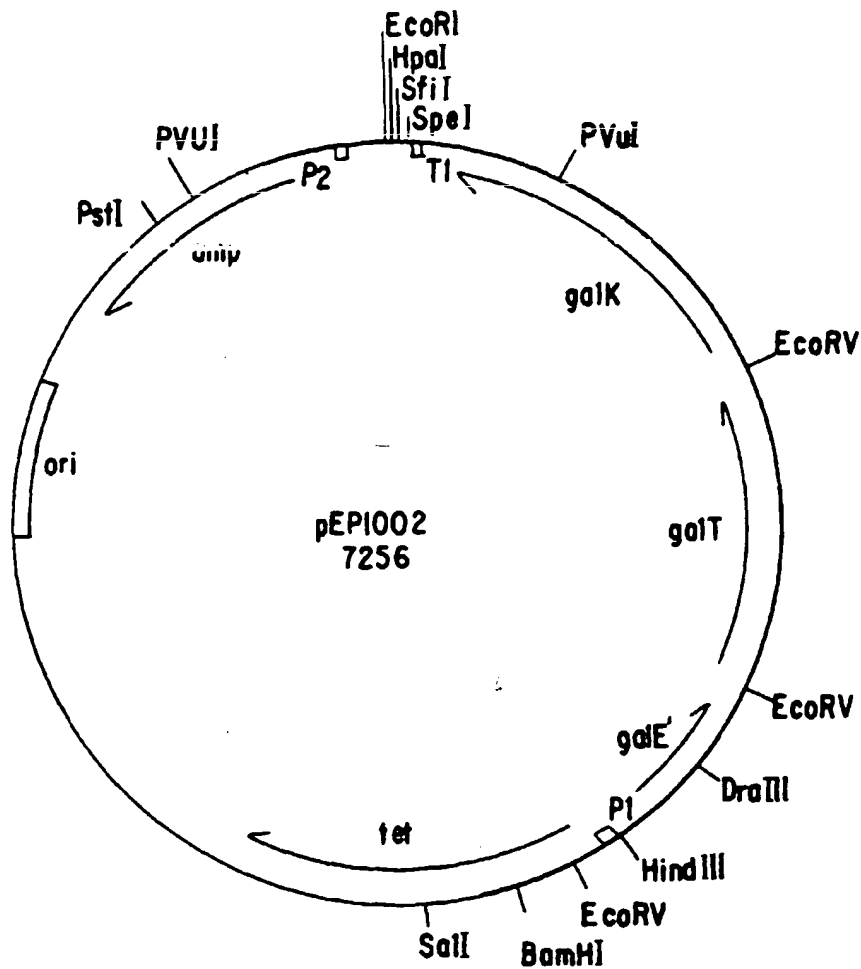
**FIG. 5.**



DELTA 4 CELLS CONTAINING pEPI001 ARE  
Amp<sup>R</sup>, Tet<sup>S</sup>, Fus<sup>R</sup>, Gal<sup>S</sup>

P1 IS pBR322 P1 IS PROMOTER;  
P2 IS Pamp

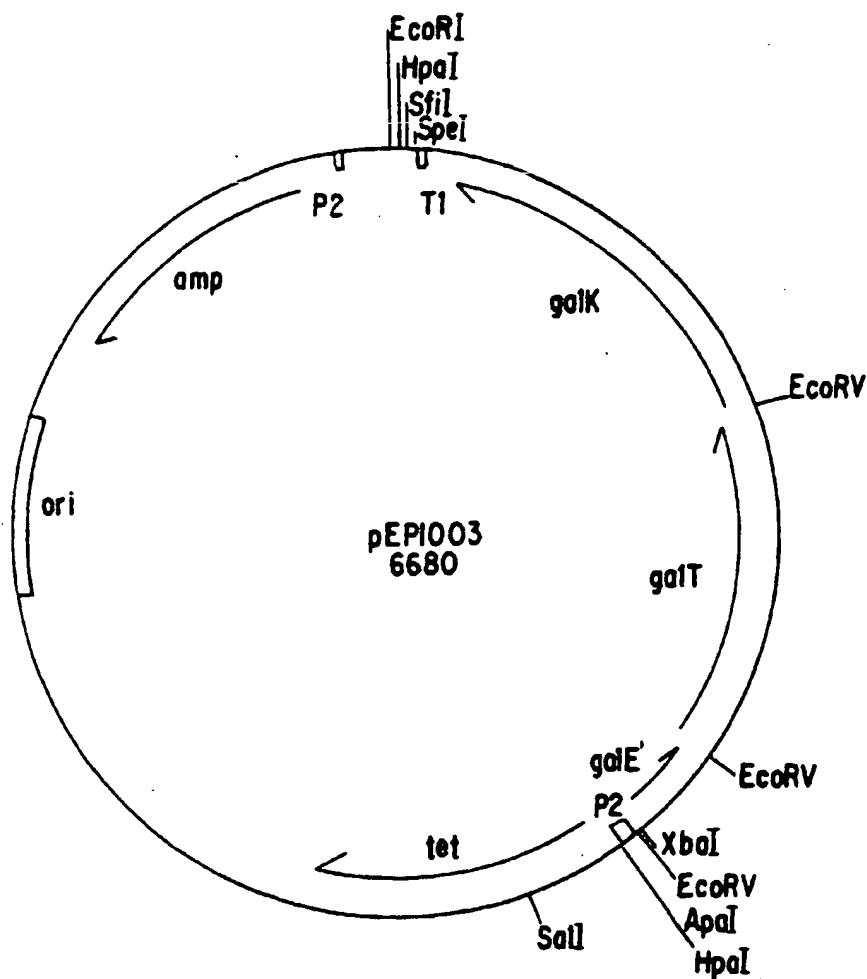
**FIG. 6.**



DELTA 4 CELLS CONTAINING pEPI002 ARE  
Amp<sup>R</sup>, Tet<sup>S</sup>, Fus<sup>R</sup>, AND Gal<sup>S</sup>

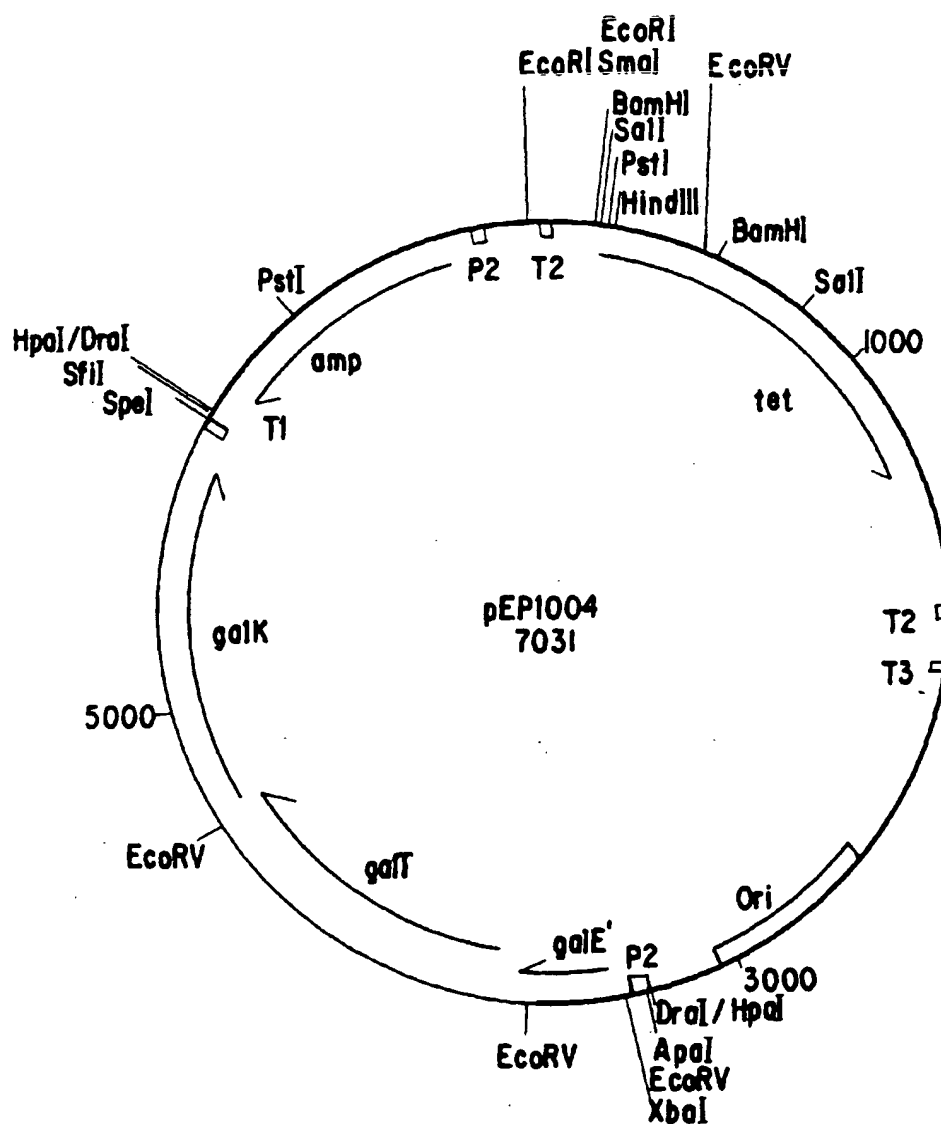
P1 IS pBR322 P1 PROMOTER; P2 IS Pamp  
T1 IS phage fd TERMINATOR; T2 IS rrnBt1;  
T3 IS rrnBt2;

**FIG. 7.**



DELTA 4 CELLS CONTAINING pEPI003 ARE  
 Amp<sup>R</sup>, Tet<sup>S</sup>, Fus<sup>R</sup>, AND Gal<sup>S</sup>  
 P2 IS Pamp  
 T1 IS phage fd TERMINATOR

**FIG. 8.**

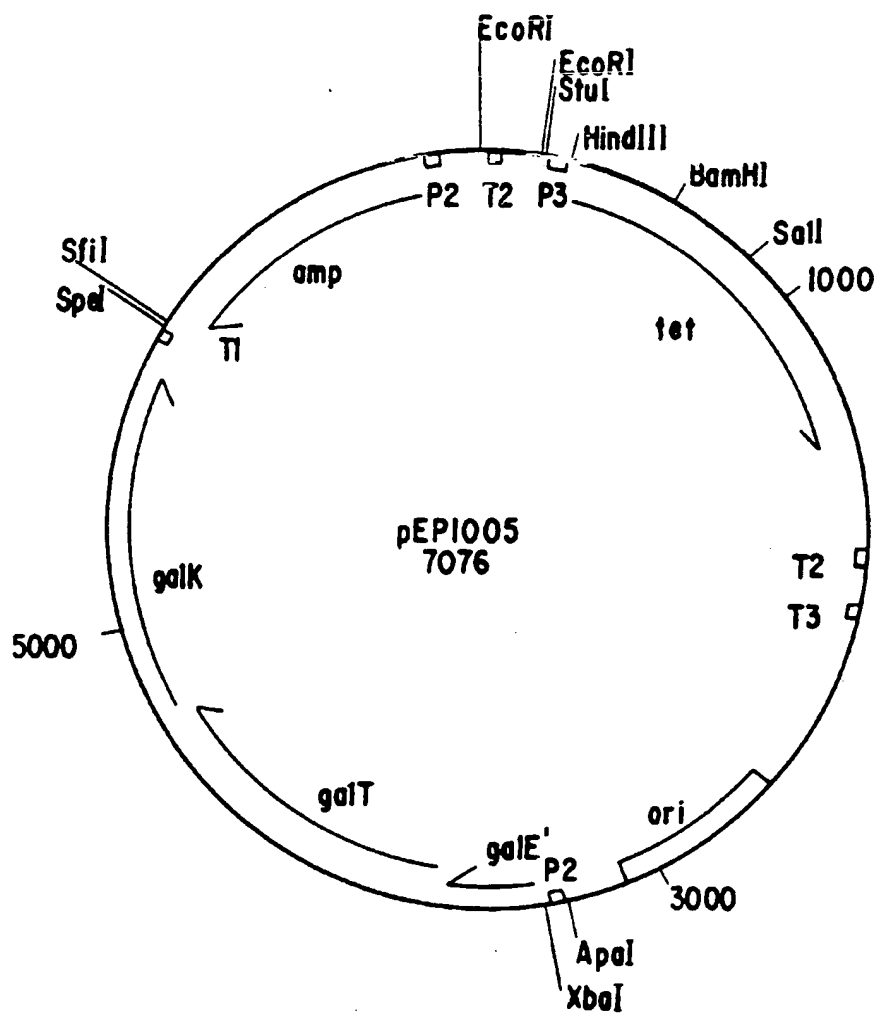
**FIG. 9.**

HB101 CELLS CONTAINING pEP1004 ARE  
Amp<sup>R</sup>, Tet<sup>S</sup>, Fus<sup>R</sup>, AND Gal<sup>+</sup>

P2 IS Pamp

T1 IS phage fd TERMINATOR, T2 IS rrnBt1, T3 IS rrnBt2,

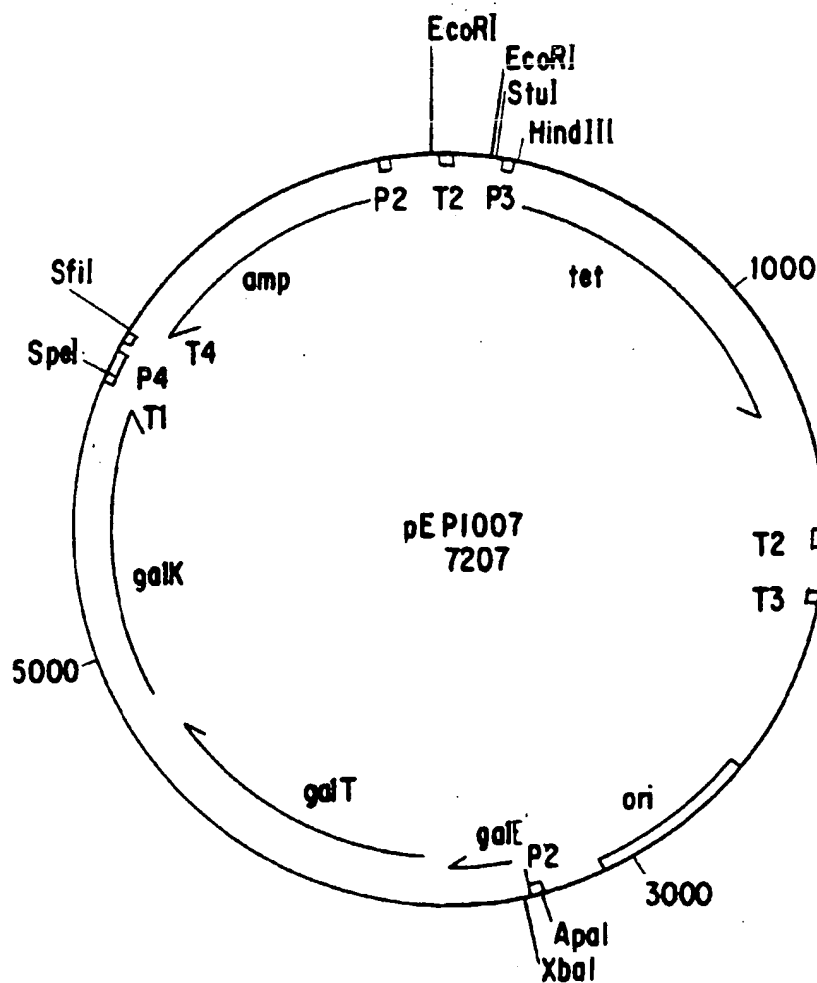


**FIG.10.**

DELTA 4 CELLS CONTAINING pEP1005 ARE  
Amp<sup>R</sup>, Tet<sup>S</sup>, Fus<sup>R</sup>, AND Gal<sup>S</sup>

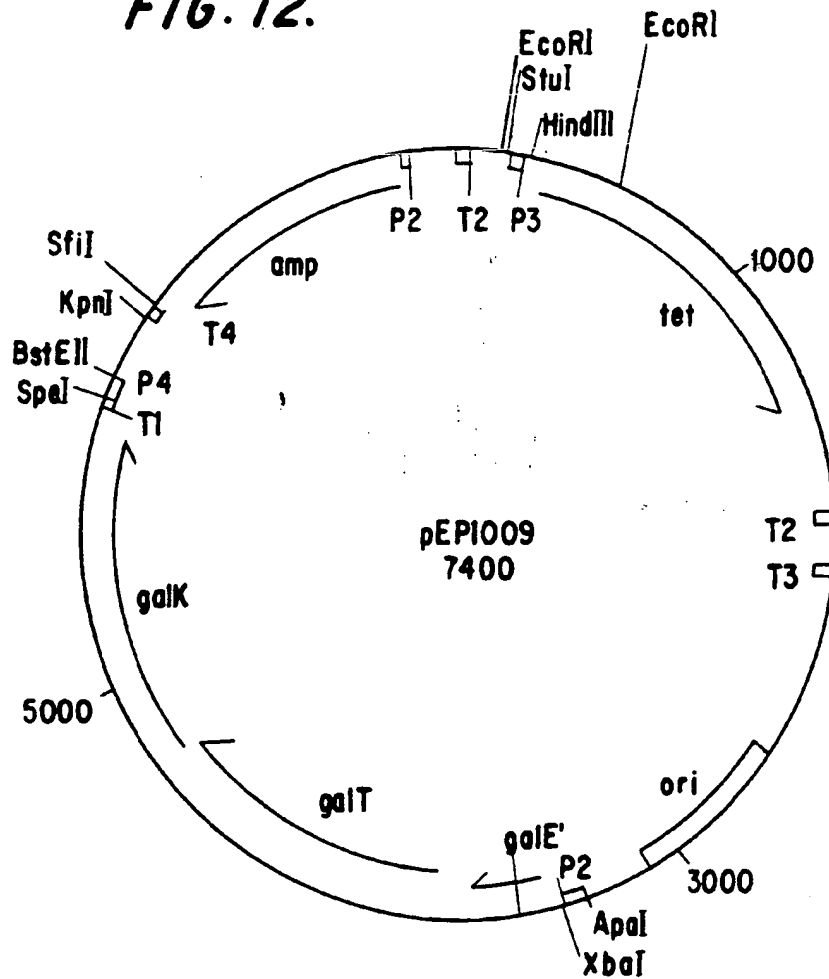
P2 IS Pamp, P3 IS Pneo

T1 IS phage fd TERMINATOR, T2 IS rrnB1 T3 IS rrnB2

**FIG. II.**

DELTA4 CELLS CONTAINING pEPI007 ARE  
Amp<sup>R</sup>, Tet<sup>R</sup>, Fus<sup>S</sup>, AND Gal<sup>S</sup>

P2 IS Pamp, P3 IS Pneo, P4 IS PlacUV5  
T1 IS phage fd TERMINATOR, T2 IS rrnB1, T3 IS rrnB2  
T4 IS trpA TERMINATOR

**FIG. 12.**

DELTA 4 CELLS CONTAINING pEPI009 ARE  
 Amp<sup>R</sup>, Tet<sup>S</sup>, Fus<sup>R</sup>, AND Gal<sup>R</sup> IN PRESENCE OF IPTG  
 Amp<sup>R</sup>, Tet<sup>R</sup>, Fus<sup>S</sup>, AND Gal<sup>S</sup> IN ABSENCE OF IPTG

P2 IS P<sub>amp</sub>, P3 IS P<sub>neo</sub>, P4 IS P<sub>lacUV5</sub>

T1 IS phage fd TERMINATOR, T2 IS rrnBt1, T3 IS rrnBt2

T3 IS trp A TERMINATOR